

000	050
001	051
002	052
003	053
004	054
005	055
006	056
007	057
008	058
009	059
010	060
011	061
012	062
013	063
014	064
015	065
016	066
017	067
018	068
019	069
020	070
021	071
022	072
023	073
024	074
025	075
026	076
027	077
028	078
029	079
030	080
031	081
032	082
033	083
034	084
035	085
036	086
037	087
038	088
039	089
040	090
041	091
042	092
043	093
044	094
045	095
046	096
047	097
048	098
049	099

CMPUT 497: Assignment 2 Report

Shouyang Zhou

University of Alberta

Edmonton, Alberta, Canada

shouyang@ualberta.ca

1 Implementation Details

The script contains three main sections, a section for constants, a class titled “N_model”, and a small set of functions used to bulk evaluate files.

The language model is implemented through the “N_model” class. The class is initialized from a filepath containing the training data. Upon initialization, the class processes the training text into a frequency distribution of n-grams and interpolation coefficients (if needed). Once the model is initialized, it can be used to evaluate a text using the “eval_text(text)” method.

Internally, the model references three frequency distributions to produce probabilities. In general, the attribute “freq_dist” is referenced to obtain counts of the actual n-gram, and the “prev_dist” is used to obtain the count of the prefix (n-1 grams) of the n-gram. These two values are then used to calculate a conditional log probability for the n-gram. Additionally, there is a unigram distribution to act as a term list. The interpolation model uses a differing set of frequency distributions titled, “freqs”, this was a later addition for the same purposes.

To evaluate a text, the model splits the evaluation text into n-grams according the model and evaluates each n-gram summing their log_2 probabilities. The sum of log-probs. are then transformed into a perplexity score. The evaluation text is pre-processed such that characters (unigrams) not found in the unigram distribution are replaced with the unknown token.

The functions provide a small evaluation framework. The function generate_evaluations_file, generates a list of models and evaluates each evaluation text using each model. It will select the model with the lowest perplexity score per evaluation text.

2 Issues Encountered

Dealing with unknown characters became an issue following the guide in J&M chapter 3. J&M recommend using an unknown token and to generate possible unknown token inclusive n-grams via preprocessing the training text, substituting infrequent terms (by some user specified criterion, ie low word occurrence) with the unknown token.

I decided to skip this unknown token preprocessing step. My rationale is that this preserves more of the actual training text (which is quite short) used to generate n-gram counts. Also, choosing the right criterion poses a confounding challenge, examining and tuning a threshold per language is too time intensive and may introduce too much bias. Using one general threshold may be too general to consider the unique characteristics in each language.

To deal with unknown terms, I will limit my unsmoothed model to an n value of 1 and seed any unigram models with a single occurrence of the unknown token.

3 Parameter Tuning

N values were tuned on the dev suffixed dataset. Iteratively, each model variant with increasing n from 1-4 was ran to produce the evaluation report upon this data. The accuracy of the model and relative perplexity scores where considered in tuning. My tests results are contained in the Excel

100	file. I aim to achieve a reasonable accuracy rate	150
101	(95%) while minimizing perplexity scores. I suspect	151
102	marginal increases in accuracy at substantial	152
103	costs in perplexity will be an indicator of data	153
104	sparseness and potential overfitting.	154
105		155
106	Mentioned earlier, the unsmoothed model has a	156
107	fixed n value of 1. This model correctly identifies	157
108	the 49 / 55 dev files and has a mean perplexity of	158
109	17.5.	159
110		160
111	Increases in N of Laplace models increased accu-	161
112	racy at the cost of “perplexity bloom”. The mod-	162
113	el’s accuracy increase (49/55, 53/55, 53/55, 54/55)	163
114	however their mean assigned perplexity scores in-	164
115	crease substantially (17.5, 16.75, 117, 426). N = 2	165
116	meets my selection criterion.	166
117		167
118	Increase of N of interpolation models were con-	168
119	sidered on the range of 2-4 given the nature of the	169
120	model. The accuracy of the model increased over	170
121	N (50/55, 52/55, 53/55), and the perplexity re-	171
122	mained approximately constant (bounded between	172
123	10-20). Examining the weights placed by the de-	173
124	leted interpolation process, all models assign uni-	174
125	grams a weighting of approximately 90% with	175
126	following ngrams being exponentially weighted	176
127	less. I choose an N value of 3.	177
128		178
129	4 Evaluation & Error Analysis	179
130	Common errors for all models involve distinctions	180
131	between dialects, for example deu_1996 vs	181
132	deu_1991 and hat_kreyol vs hat_popular. Examin-	182
133	ing these two as case studies, both training da-	183
134	tsets share very similar character sets and words.	184
135	Notably, the deu training sets differ by the re-	185
136	placement of a single character. In these cases, the	186
137	small variations in the training datasets, notably	187
138	length, may favor certain models.	188
139		189
140	Errors were correlated in general; some dev files	190
141	appear be simply harder to distinguish overall. Er-	191
142	rors within the same model class, across n are	192
143	most highly correlated.	193
144		194
145	Harder files are often texts that draw similarities	195
146	to other languages. For example, “udhr-	196
147	sco.txt.dev” appears to draw on many English like	197
148	words. I suspect one confounding factor to be text	198
149	forms that in some form are a spelling simplifica-	199
	tion of another text form. If this simplification fa-	