

1 Tokenizing

The sonnets were initially tokenized by words using the NLTK's function `word.tokenizer` because word is the smallest entity that we can tokenize in the sonnets such that we can still regenerate sentences with the tokens later. We created two types of data set. First, we consider each line as a sequence, and second, we consider each sonnet as a sequence.

The first change we made after running the learning algorithm was to remove all punctuations. We decided to generate the sentences backward using a seed (i.e. a word) to achieve accurate ending rhythm. Thus, the punctuations were mostly irrelevant because they showed up mostly in the end. In addition, we also find that the model with punctuations generated sentences that sometimes had brackets that did not match or commas that did not make gramatical sense. Therefore, punctuations were removed in our final version.

Furthermore, we counted the bigrams within the sonnets and replaced each pair of consecutive words that showed up frequently in the sonnets as one "word". These word pairs were added into the list of all word tokens. In total, 29 pairs of words were added. Each pair appeared in the sonnets at least 15 times. The top five examples were "my love", "thou art", "my self", "in the ", "that i" which appeared 41, 33, 29, 29, and 28 times respectively.

Lastly, we created a rhythm dictionary that was used in poem generation later to create sentences that had the correct rhymes. For the rhythm dictionary, we collected all ending words in each line of the sonnets that rhymed. For example, "increase" and "decrease" is a pair of words that rhymed. So, they are stored as a tuple in the rhythm dictionary.

2 Algorithm

For unsupervised training of our Hidden Markov Model, we used the Baum-Welch algorithm. We used the TA solution from HW #5 for the implementation. The main three parameters we control were (1) the number of EM iterations, (2) the number of hidden states to use, and (3) whether to train on lines of each poem or the poems themselves. Regarding (1), we used 1000 iterations, because this was enough iterations to show convergence to a very reasonable precision. Regarding (2) and (3), we trained several different HMMs, each with a different number of hidden states (8, 10, 12, 15, 16, 20, 25, 30). For each choice of hidden states, we also trained on each line of the poem, and on each poem itself. We then visualized the HMMs through both the transition matrix and observation matrix (discussed in section 5), and examined the generated poems to determine the optimal number of hidden states based on a qualitative analysis of the HMM structure, and whether to train on lines or poems.

We eventually chose an HMM with 20 hidden states trained on the individual lines of the poem, based on the underlying structure of the HMM and the qualitative quality of the generated poems. When we assess the underlying structure of the HMM, we refer to meaningful patterns in the transition matrix and observation matrix (also whether states generate words with discernible patterns). This underlying structure is discussed in section 5. When using fewer or more states, this structure was less clear, resulting in our decision to use 20 hidden states.

When generating our poem, we also needed to keep track of both syllable count. To do this, we used the NLTK CMU dictionary to keep syllable counts.

3 Poem Generation

First Bare-Bones Implementation

For our initial bare-bones implementation, we originally generated poems by the following steps:

1. Pick a uniformly random states to begin at
2. Initialize an empty line and initialize the syllable count of this line to zero.
3. Sample the next state based on the previous state and the transition matrix probabilities.
4. Add a word to the line based on the current state and the observation matrix probabilities
5. Count the number of syllables in the added word and add it to a syllable count
6. Repeat steps 3-5 until we have created a line with at least 10 syllables. Save this line to our poem.
7. Repeat steps 1-6 until we have generated 14 lines for our poem.

To get the syllable count in step 5, we used the NLTK CMU dictionary to keep syllable counts. By following the steps above, we were able to generate 14-line poems with each line containing *at least* 10 syllables.

4 Poem Generation/Additional Goals

However, we wanted to generate poems with *exactly* 10 syllables. To do this we modified step 6 to get the following algorithm:

Implementation with Syllable Counts (No Rhyme)

However, we wanted to generate poems with *exactly* 10 syllables. To do this we modified step 6, and get the following algorithm:

1. Pick a uniformly random states to begin at
2. Initialize an empty line and initialize the syllable count of this line to zero.
3. Sample the next state based on the previous state and the transition matrix probabilities.
4. Add a word to the line based on the current state and the observation matrix probabilities
5. Count the number of syllables in the added word and add it to a syllable count
6. Repeat steps 3-5 until we have created a line with at least 10 syllables.
 - (a) If we have exactly 10 syllables. Add the line to the poem
 - (b) If we have >10 syllables, remove the last word, go to the previous state and repeat steps 3-5.
7. Repeat steps 1-6 until we have generated 14 lines for our poem.

This gave poems with exactly 10 syllables. The tradeoff here was that the last word of the line was often restricted in its number of syllables (i.e. if we had 8 syllables in our line already, we had to pick a last word with less than 3 syllables). Thus the last word was picked from a probability distribution that did not exactly match that given by the bare-bones HMM.

Implementation with Syllable Count and Rhyme

Finally, we want to implement rhyme. To do this, we had to generate our poems backward using the following algorithm:

1. Create a dictionary of rhyming pairs based on the rhyming pairs in the shakespeare poems we trained on.
2. For lines 1, 2, 5, 6, 9, 10, 13, randomly choose a word from the rhyming dictionary. For all other lines, choose the rhyming word that goes with the corresponding line (e.g. for line 4, choose the rhyming word from our rhyming dictionary that is paired with the word chosen for line 2).
3. Sample the state based on calculated probability $P(y \mid x)$ where y represents a hidden state and x represents our observation of the first word. We can calculate $P(y \mid x)$ using the formula $P(y^1 = z \mid x) = \frac{\alpha_z(1)\beta_z(1)}{\sum_{z'} \alpha_{z'}(1)\beta_{z'}(1)}$. We use the *forward*, *backward* functions we wrote for HW #5 in order to get α and β . Once we have our chosen first word and have calculated $P(y \mid x)$, we can sample from this distribution to get y^1 .
4. Add a word to the line based on the current state and the observation matrix probabilities
5. Count the number of syllables in the added word and add it to a syllable count
6. Sample the next state based on the previous state and the transition matrix probabilities.
7. Add a word to the line based on the current state and the observation matrix probabilities
8. Count the number of syllables in the added word and add it to a syllable count
9. Repeat steps 6-8 until we have created a line with at least 10 syllables.
 - (a) If we have exactly 10 syllables. Add the line to the poem
 - (b) If we have >10 syllables, remove the last word, go to the previous state and repeat steps 6-8
10. Reverse the generated line (because we have generated the line backwards to ensure proper rhyming), and add it to the poem.
11. Repeat steps 2-10 until we have generated 14 lines.

Generation Result/Discussion

Using this algorithm, we were able to successfully automatically generate poems using an HMM trained on several shakespeare poems that obeyed the syllable count and rhyming pattern of a shakespearean sonnet. However, by implementing rhyme in the above fashion, we had to restrict the choice of first word in every line to one from our rhyming dictionary. Thus the structure of our generation algorithm using HMM was different. In our first step, rather than transitioning to an initial state (based on A) and generating a word (based on O), we had to first choose the word, and then sample the most likely state.

Generating the poems in this fashion, enforcing both rhyming scheme and syllable count, had the effect:

1. Since we chose to cut off the line at 10 syllables as we did, the last word of the line was often restricted in its number of syllables (i.e. if we had 8 syllables in our line already, we had to pick a last word with less than 3 syllables). Thus the last word was picked from a probability distribution that did not exactly match that given by the bare-bones HMM, in order to ensure each line had 10 syllables.

2. Because we were considering transitions backward (to the previous word in the line), the state transition matrix became less intuitive to follow since reading sentences backwards results in different transition structure and is less natural to humans.
3. Sometimes the first word of each line did not fit as a first word, since it was the last word generated in our backwards generation scheme.
4. Our first word of each line was limited to our rhyming dictionary. Furthermore, we had to calculate $P(y | x)$ in order to get the first state from the first word,

Below is one of the poems we generated with the proper rhyming scheme and syllable count (chosen for Piazza submission):

Love my self alteration side granting
 Plague for and bright oaths but sit faces moan
 Sky grow'st best wilt or not to on wanting
 Up turns abide votary shall upon

Not bide canst the that which to oppressed
 Am too black try is fair am be increase
 Usurer in chest state party age rest
 Change of in not to knows my muse decease

Can canker be war him forth him power
 Days loss is ward clock dead eloquence pride
 I still verse with abide play'st call flower
 Again cruel told kindness odour one side

Than to that i west slow and sounds a reap
 All friend great poor tear that all that i leap

10

5 Visualization and Interpretation

The top 10 words that associate with the hidden state are given in Table ?? (the numbers below the words denote the probability). The number in bracket is the total number of top words to form approximately 0.5 probability. Hidden state 11 and 12 have probabilities that are more diffused and they includes more words that are nouns, adverbs or adjectives. Hidden state 1, 8, 9, and 18 have a more concentrated probabilities, and they tend to be words that are used to connect phrases like particles, conjunctions, and pronouns. Hidden state 10 includes the 5 “wh” words for questions.

Figure ?? shows the observation matrix in which each row is normalized to 0 and 1 for clearer color contrast. Each column represents a token and the columns are sorted by the frequency of the tokens. The top plot shows all tokens, and the bottom plot shows the first 100 tokens. From the plots, we can see that the most common words in general have a higher probability than the less common words. Words like “and” which are very frequently used (473 times) in the sonnets show up in a few hidden states with high probability.

If we also consider the transition matrix (Figure ??), we can observe some interesting state transitions that makes grammatical sense. For example, hidden state 11 is most likely to transition to hidden state 19 with probability 0.9330. Combinations of the top 4 words “of”, “in”, “to”, and “not” in hidden state 19 with

most of the top 10 words in hidden state 11 (e.g. love, i, one, hath, beauty, this, and most) seem to make grammatical sense. Another example is hidden state 1 is most likely to transition to hidden state 12 with probability 0.9836. Combinations of the top words “is” and “time” in hidden state 12 with “the”, “my”, “a”, “so”, “such”, “their”, “his”, and “of the” in hidden state 1 also make grammatical sense.

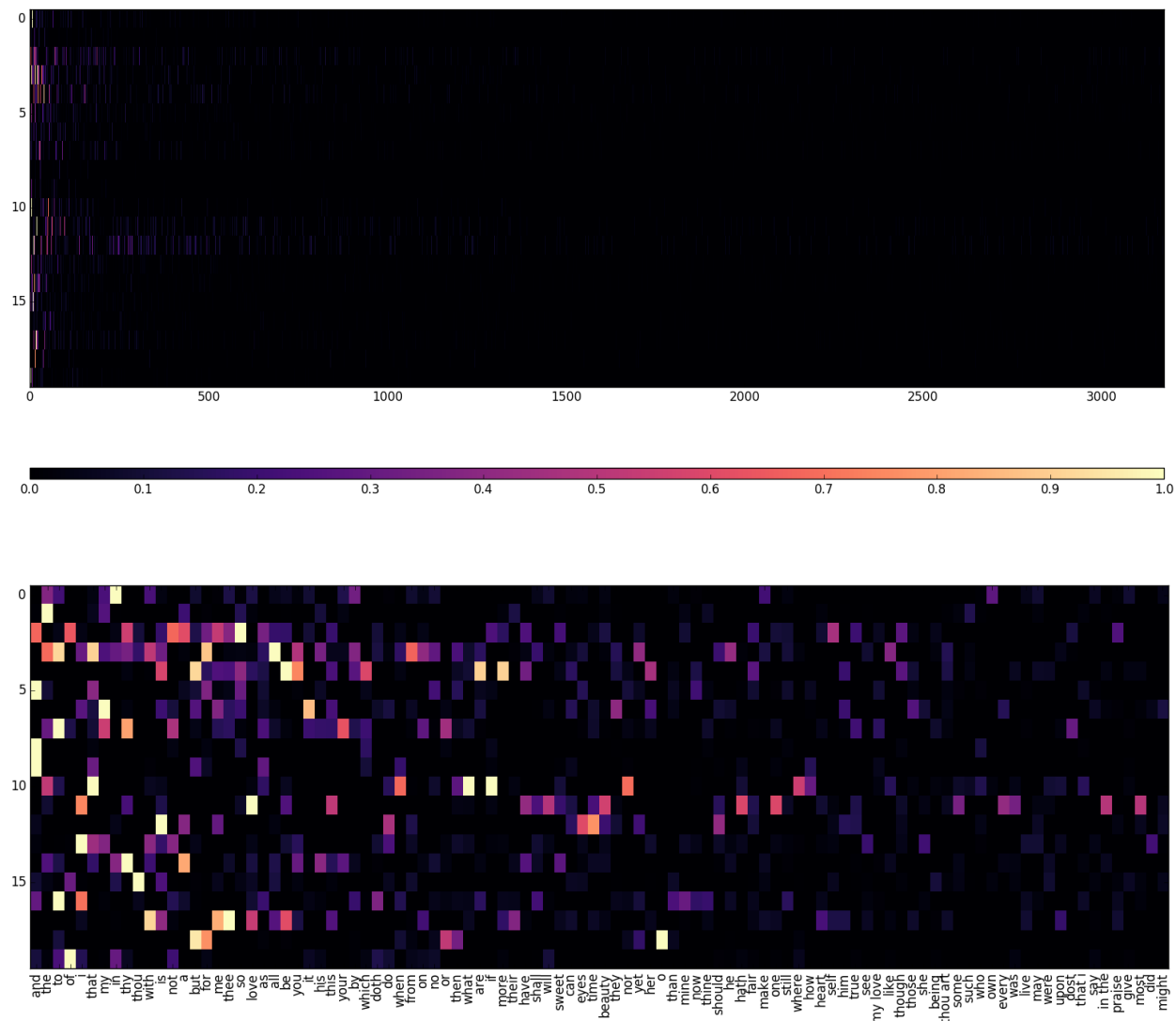


Figure 1: Normalized observation matrix. The top plot shows all tokens, and the bottom plot zooms into the first 100 tokens.

Table 1: Top 10 words for hidden states and the probability. The number in bracket is the total number of top words to form approximately 0.5 probability.

State	Top 10 Words									
0 (26)	in 0.1054	the 0.0399	by 0.0352	own 0.0300	my 0.0250	with 0.0244	make 0.0230	to 0.0230	show 0.0149	your 0.0140
1 (9)	the 0.2325	my 0.0604	a 0.0481	so 0.0337	such 0.0319	their 0.0283	his 0.0251	of the 0.0204	well 0.0164	thine 0.0158
2 (50)	so 0.0364	not 0.0248	and 0.0247	of 0.0234	a 0.0223	thy 0.0203	me 0.0195	self 0.0169	youth 0.0143	thee 0.0140
3 (23)	all 0.0454	to 0.0417	that 0.0413	for 0.0400	the 0.0313	from 0.0305	with 0.0236	you 0.0221	he 0.0192	yet 0.0187
4 (30)	be 0.0444	but 0.0407	more 0.0402	are 0.0401	you 0.0332	which 0.0262	is 0.0257	her 0.0205	so 0.0179	i am 0.0177
5 (35)	and 0.0985	that 0.0387	for 0.0361	so 0.0311	no 0.0217	now 0.0210	then 0.0152	eyes 0.0150	to me 0.0148	when i 0.0128
6 (15)	my 0.0993	it 0.0854	they 0.0403	me 0.0371	i 0.0271	but 0.0259	those 0.0257	this 0.0228	her 0.0226	is 0.0205
7 (21)	to 0.0702	thy 0.0565	your 0.0442	my 0.0384	not 0.0333	or 0.0318	the 0.0217	dost 0.0215	on 0.0201	as 0.0167
8 (2)	and 0.3856	which 0.0657	so 0.0509	who 0.0439	therefore 0.0202	if 0.0189	even 0.0136	not 0.0134	or 0.0112	for 0.0106
9 (3)	and 0.2786	that 0.0876	as 0.0825	but 0.0707	when 0.0523	how 0.0423	which 0.0287	since 0.0221	for 0.0193	can 0.0151
10 (10)	if 0.0723	what 0.0719	that 0.0719	nor 0.0502	when 0.0492	where 0.0400	the 0.0380	let 0.0336	why 0.0311	how 0.0221
11 (61)	love 0.0354	i 0.0265	one 0.0209	hath 0.0208	beauty 0.0200	this 0.0185	most 0.0178	in the 0.0177	will 0.0175	part 0.0149
12 (81)	is 0.0305	time 0.0235	eyes 0.0183	do 0.0143	should 0.0126	a 0.0116	sun 0.0105	summer 0.0101	age 0.0100	best 0.0095
13 (31)	i 0.0914	that 0.0403	with 0.0359	not 0.0352	my 0.0342	do 0.0220	did 0.0197	she 0.0193	see 0.0180	is 0.0168
14 (18)	thy 0.0958	a 0.0783	his 0.0379	in 0.0368	you 0.0294	your 0.0269	sweet 0.0267	have 0.0252	the 0.0237	with 0.0196
15 (31)	thou 0.1502	of 0.0474	is 0.0367	in 0.0202	all 0.0170	being 0.0162	and 0.0136	were 0.0133	from 0.0126	day 0.0114
16 (19)	to 0.1011	i 0.0712	doth 0.0368	mine 0.0312	and 0.0304	than 0.0244	shall 0.0223	as 0.0199	by 0.0197	thine 0.0194
17 (28)	thee 0.0661	with 0.0605	me 0.0547	be 0.0405	love 0.0358	is 0.0269	their 0.0249	all 0.0176	this 0.0168	heart 0.0155
18 (3)	o 0.1761	but 0.1683	for 0.1355	or 0.0924	then 0.0595	yet 0.0321	against 0.0303	save 0.0241	to 0.0198	when 0.0169
19 (14)	of 0.1831	in 0.0562	to 0.0431	not 0.0380	i 0.0269	doth 0.0249	and 0.0205	will 0.0200	do 0.0166	of thy 0.0150

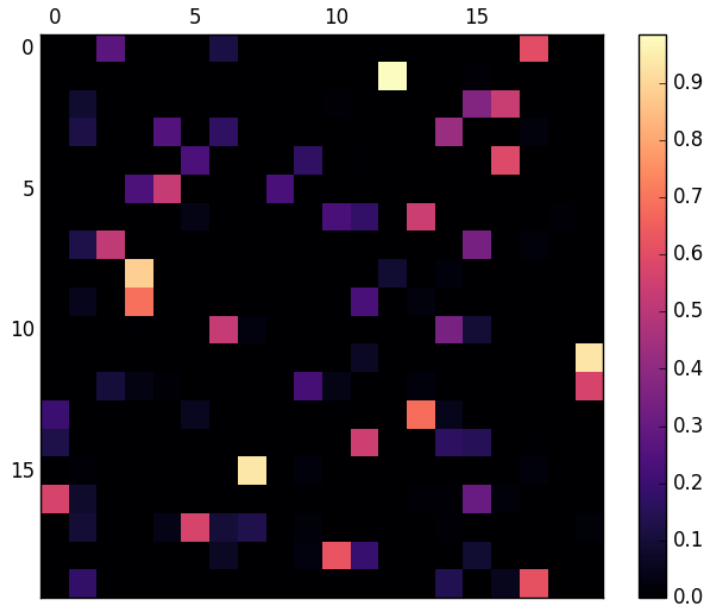


Figure 2: State transition matrix

6 Conclusion

The poems we generated had some structure and seemed to generate some grammatically sensical lines. It was usually able to string together short sequences of words in a grammatically coherent way (though this often broke down with longer sequences). In the “Poetry Generation/Additional Goals” section, we discussed the trade-offs we had to make to enforce rhyme and syllable count in our poem. A few other shortcomings of our poems, and potential remedies, are discussed below.

1. Because we trained our HMM on individual lines, the generated poems had no theme (beyond the coincidental). Thus each line often showed little thematic correlation to adjacent lines. This could be remedied by training on poems rather than lines, but we ended up choosing to train on lines, because the grammatical structure of generated lines was more coherent. With more data, we could consider training on full poems though.
2. By nature of generation with HMMs, words can only have direct relation with the previous word. Therefore it becomes likely for long lines to wander, and difficult to string together many words for a single coherent idea. Thus, in our poems, seldom are there long strings of words forming a single coherent idea.
3. Some words or pairs of words were extremely frequent in the poems, and the prominence of common words like ‘the’, ‘be’, and ‘he’ often threw off the structure of sentences and resulted in part-of-speech sequences that did not make great sense. To remedy this, we could try training our model to learn part-of-speech tags. Google’s Speech API is able to identify part-of-speech tags in a given sentence, so we could use this to help our model understand words’ part-of-speech tags and generate sentences with proper part-of-speech sequences (e.g. pronoun, noun, adverb, verb, noun).

Despite these issues, we were happy that our poems showed some likeness to shakespearean poetry with mostly sensical lines. Furthermore, we were able to visualize our Hidden Markov Model and extract some

structure and learn some interesting things about shakespeare's style and word structure (e.g. he really enjoys "my love").

We broke up the work...