

1 Overview

- In this project, we will attempt to generate Shakespearean sonnets by training a HMM on the entire corpus of Shakespeare's sonnets.
- This miniproject is due 5pm, March 10th, via Moodle.
- You can work in groups of up to three. We encourage you to use the [Search for Teamates feature on Piazza](#) to help you find teammates. You may keep the same group as in mini-project 1.
- You may use any language you want, but you must submit a report including documented code that lays out everything you did. We recommend Python for this assignment.
- You are required to share one poem with the class on Piazza. There will be a small competition for who can generate the best poem!

2 Introduction

William Shakespeare is perhaps the most famous poet and playwright of all time. Shakespeare is known for works such as Hamlet and his 154 sonnets, of which the most famous begins:

*Shall I compare thee to a summer's day?
Thou art more lovely and more temperate:*

Shakespeare's poems are nice for the purpose of generative modeling because they follow a very specific format, known as the Shakespearean (or English) sonnet¹. Each sonnet is 14 lines, spread into 3 quatrains (section with 4 lines) followed a couplet (section with 2 lines). The third quatrain is generally known as the *volta*² and has a change in tone or content. Shakespearean sonnets have a particular rhyme scheme, which is *abab cdcd efef gg*.

Shakespearean sonnets also follows a specific meter called *iambic pentameter*³. All lines are exactly 10 syllables long, and have a pattern of unstressed stress. For example, the famous Sonnet 22 begins:

Stress	x	\	x	\	x	\	x	\	x	\
Syllable	Shall	I	com -	pare	thee	to	a	sum-	mer's	day?

Here, each x represents an unstressed syllable and every \ represents a stressed syllable. Try saying it out loud!

The goal for this project is to generate poems that Shakespeare may have written by training a HMM on his 154 sonnets. His sonnets are available in the data file `shakespeare.txt`.

¹https://en.wikipedia.org/wiki/Sonnet#English_.28Shakespearean.29_sonnet

²https://en.wikipedia.org/wiki/Volta_%28literature%29

³https://en.wikipedia.org/wiki/Iambic_pentameter

3 Unsupervised Learning (30 points)

The first task is to implement the unsupervised learning algorithm for HMMs. Use the forward-backwards algorithm to train a HMM on the dataset. To determine the number of hidden states, you should try training models with different numbers of hidden states, and see how they perform. You may base the code from the homework solution (which only solves the supervised M-step).

You should preprocess the dataset before you train on it. How you preprocess is completely up to you. Here are a couple of questions to help you decide how to do preprocessing: How will you tokenize the data set? What will consist of a singular sequence, a poem, a stanza, or a line? Do you keep some words tokenized as bigrams? Do you split hyphenated words? How will you handle punctuation?

Report:

Your report should contain a section dedicated to preprocessing. Explain your choices, as well why you chose these choices initially. What was your final preprocessing? What changed as you continued on your project? What did you try that didn't work? Also write about any analysis you did on the dataset to help you make these decisions.

Your report should also contain a section highlighting your unsupervised learning. How did you choose the number of hidden states? How did you tokenize your words, and split up the data into separate sequences?

NOTE: You may use an off-the-shelf implementation of HMMs for the other parts of the miniproject if you cannot get your implementation working properly. However, you must write your own code for this part of the project.

4 Visualization & Interpretation (30 points)

The next section of this project deals with interpreting and visualizing the learned model. Our goal is to interpret what the hidden states and transitions capture about the data. Use the learned observation matrix and transition matrix to determine what words associate most with each hidden state, and how the hidden states interact with each other. Do the hidden states represent parts of speech, stressed or unstressed words, number of syllables? What about anything else you can think of? You may use any open source tools to help you perform some of the analysis. You can get syllable counts and syllable stress information from CMU's Pronouncing Dictionary available on [NLTK](#)

Report:

In your report, you should explain your interpretation of how a Hidden Markov Model learns patterns in Shakespeare's texts. You should briefly elaborate on the methods you used to analyze the model. In addition, for at least 10 hidden states give a list of the top 10 words that associate with this hidden state and state any common features these groups. Furthermore, try to interpret and visualize the learned transitions between states. A possible suggestion is to draw a transition diagram of your markov model and give descriptive names to the states. Feel free to be creative with your visualizations, but remember that accurately representing data is still your primary objective. Your figures, tables, and diagrams should contribute to a discussion about your model.

5 Poetry Generation (20 points)

Some theory

Remember that the core of a HMM is the transition matrix and the observation matrix. Given a current state y_0 , we can generate the next state by randomly choosing a state from the row in the transition matrix based on the probability of transitioning to that state. Now, with the next state, we can generate a word by choosing randomly based on each word's probability of generated from that state.

Naive Poem Generation from HMM's

Write a program that generates a 14-line sonnet from a HMM model. You will need to choose one sonnet that you generate and share it with the rest of the class on Piazza under the tag `project2`. The TAs will read over your submissions and choose the best computer generated sonnet. Note that the poem that you submit for the competition does not need to be from the naive poem generation, and can be from a later improved HMM model. However, the poem you submit must be computer generated. You may update your poems until the deadline for the project.

Report:

In your report, describe your algorithm for generating the 14 line sonnet. Include at least one sonnet generated from your unsupervised trained HMM in your final report as an example. You should comment on the quality of generating poems in this naive manner. How accurate is rhyme, rhythm, and syllable count to what a sonnet should be? Do your poems make any sense? Does it retain Shakespeare's original voice? How does training with different number of hidden states effect the poems generated (in a qualitative manner)? For the good qualities that you describe, also discuss how you think the HMM was able to capture these qualities.

6 Additional Goals (20 points)

"Though this be madness, yet there is method in't" - Hamlet Act 2, scene 2

So far, the poems you have generated using the naive HMM generation are probably not very good as sonnets. In this section, you will explore methods of improving your poems or extending them. **You do not need to attempt all of the tasks below for full marks on this section.** If you have ideas for other improvements to the poetry generation not listed here, feel free to talk to a TA and work on it. The sky is the limit.

Report:

Talk about the extra improvements you made to your poem generation algorithm. What was the problems you were trying to fix? How did you go about attempting to fix it? Why did you think that what you tried would work? Did your method succeed in making the sonnet more like a sonnet? If not, why do you think what you tried didn't work? What tradeoffs do you see in quality and creativity when you make these changes?

Rhyme

Introducing rhyme into your poems is not actually that difficult. Since the sonnet follows a strict rhyming patterns, we can figuring out what rhymes Shakespeare uses by looking at the last word of each pair of lines that rhyme and add this to some sort of rhyming dictionary. Then, we can generate two lines that rhyme by seeding the end of the line with words that rhyme, and then performing the HMM generation in the reverse direction.

Meter

One way to incorporate meter is by creating states that represent the stresses of a word and limiting transitions between stressed and unstressed words. For example, if a word ends in a stressed syllable, its state should not transition to a state that represents words that start with a stressed syllable. You can also guarantee a syllable count by using supervised learning and labeling words by syllable and stress, and counting syllables when generating your poem. However, you may find that a more constrained HMM may produce lower-quality sentence structure. If you use too many states, the HMM may lose variety in its generation. To find a happy medium, try semi-supervised learning.

Incorporating additional texts

A powerful feature of HMMs is the ability to combine texts from different sources, with potentially silly results.⁴ We have also provided the Amoretti⁵ by Spenser, a contemporary poet of Shakespeare. All 139 of Spenser's sonnets in the Amoretti follow the same rhyme scheme and meter as Shakespeare's sonnets.

Generating other poetic forms

It may be an endeavor to try to generate other poetic forms using your HMMs. Can you generate Haikus? How about Petrarchan sonnets, limmericks?

Choose your own!

This project is meant for you to have fun and explore new ideas. Talk to the TA's about your own ideas of how to make the poems better, and try it out. We may award bonus points for creativity.

7 Additional Resources

- [TED talk: Can a computer write poetry?](#)
- [Botpoet](#)
- [Natural Language Processing Toolbox](#)
- [Markov Constraints for Generating Lyrics with Style](#)
- [Unsupervised Rhyme Scheme Identification Hip Hop Lyrics Using Hidden Markov Models](#)

⁴[King James bible mixed with SICP](#)

⁵<https://en.wikipedia.org/wiki/Amoretti>