

(520|600).666
Information Extraction from Speech and Text

Programming Assignment # 1

Due March 7, 2024.

You will model letters of English text using hidden Markov models. Some ordinary text has been selected and, to keep matters simple, all numerals and punctuation have been purged, upper case letters have been lower-cased, and inter-word spacing, new lines and paragraph breaks, have all been normalized to single spaces. The alphabet of the resulting text is therefore the 26 lower case English letters and the white-space, and is formally denoted as $\mathcal{Y} = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, \dots, \mathbf{z}, \#\}$, with $\#$ denoting the white-space character.

The text is 35,000 characters long, and has been divided into a 30,000 character *training* set, named **A**, and a 5,000 character *test* set, named **B**.

1. Model this text with a fully connected 2-state HMM, with states ① and ②.

Let t_1 denote the transition ① \rightarrow ②, and t_2 denote the self-loop ① \rightarrow ①. Similarly, let t_3 denote the transition ② \rightarrow ①, and t_4 denote the self-loop ② \rightarrow ②, so that the transition probability matrix may be written as

$$\mathbf{p} = \begin{bmatrix} p(\textcircled{1}|\textcircled{1}) & p(\textcircled{2}|\textcircled{1}) \\ p(\textcircled{1}|\textcircled{2}) & p(\textcircled{2}|\textcircled{2}) \end{bmatrix} \equiv \begin{bmatrix} p(t_2) & p(t_1) \\ p(t_3) & p(t_4) \end{bmatrix}$$

Let the initial state of the Markov chain be either ① or ② with *equal probability*.

Let the emission probabilities be associated with the states, i.e. let

$$q(y|t_2) \equiv q(y|t_3) \equiv q(y|\textcircled{1}) \quad \text{and} \quad q(y|t_1) \equiv q(y|t_4) \equiv q(y|\textcircled{2}).$$

Use the Baum-Welch algorithm and the training text **A** to estimate the probabilities $p(t_j)$, $j = 1, 2, 3, 4$, and the emission probabilities $q(y|s)$, $y \in \mathcal{Y}$ and $s \in \{\textcircled{1}, \textcircled{2}\}$.

- (a) Initialize the transition probabilities to be *slightly different from uniform*, as

$$p(t_1) = 0.51 = p(t_3) \quad \text{and} \quad p(t_2) = 0.49 = p(t_4).$$

Initialize the emission probabilities to also be *slightly different from uniform*, as

$$q(\mathbf{a}|\textcircled{1}) = q(\mathbf{b}|\textcircled{1}) = \dots = q(\mathbf{m}|\textcircled{1}) = 0.0370 = q(\mathbf{n}|\textcircled{2}) = q(\mathbf{o}|\textcircled{2}) = \dots = q(\mathbf{z}|\textcircled{2}),$$

$$q(\mathbf{a}|\textcircled{2}) = q(\mathbf{b}|\textcircled{2}) = \dots = q(\mathbf{m}|\textcircled{2}) = 0.0371 = q(\mathbf{n}|\textcircled{1}) = q(\mathbf{o}|\textcircled{1}) = \dots = q(\mathbf{z}|\textcircled{1}),$$

$$q(\#|\textcircled{1}) = 0.0367 = q(\#|\textcircled{2}).$$

What would happen if all probabilities were set to be uniform, i.e. $\frac{1}{2}$ and $\frac{1}{27}$?

- (b) Plot the average log-probability of the training and test data after k iterations,

$$\frac{1}{|\mathbf{A}|} \log P_k(\mathbf{A}) \quad \text{and} \quad \frac{1}{|\mathbf{B}|} \log P_k(\mathbf{B}),$$

as a function of the number of iterations, for $k = 1, 2, \dots, 600$.

- (c) Plot the emission probabilities of a few particular letters for each state, e.g.

$$q_k(\mathbf{a}|\textcircled{1}) \text{ versus } q_k(\mathbf{a}|\textcircled{2}) \quad \text{and} \quad q_k(\mathbf{n}|\textcircled{1}) \text{ versus } q_k(\mathbf{n}|\textcircled{2}),$$

as a function of the number of iterations, for $k = 1, 2, \dots, 600$.

- (d) Study the emission probability distributions $q_{600}(\cdot|\textcircled{1})$ and $q_{600}(\cdot|\textcircled{2})$ to see where they differ the most, as well as how the transition probabilities differ from their initial values. Try to explain what the machine has learned about English text.

2. *Increasing Model Complexity:* Repeat the Exercises 1(a) through 1(d) with a fully connected 4-state HMM. Modify the initialization in 1(a) to account for 4 states.
3. *Alternate Initialization of Output Probabilities:* HMM estimation is sometimes sensitive to the initialization of the model parameters. You will now investigate an alternative to the initialization of Exercise 1(a).

- (a) Compute the relative frequency $q(y)$ of the letters in \mathcal{Y} from the entire text \mathbf{A} .
- (b) Generate a vector of *random* numbers $r(y)$, compute the average $\bar{r} = \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} r(y)$, and use it to create a zero-mean *perturbation* vector $\delta(y) = r(y) - \bar{r}$.
- (c) Choose a small $\lambda > 0$, though not too small, such that both

$$q(y|\textcircled{1}) = q(y) - \lambda\delta(y) > 0 \quad \text{and} \quad q(y|\textcircled{2}) = q(y) + \lambda\delta(y) > 0 \quad \forall y \in \mathcal{Y}.$$

Note: $q(\cdot|\textcircled{1})$ and $q(\cdot|\textcircled{2})$ are bona fide probability assignments on \mathcal{Y} . (Why?)

Use the two $q(y|s)$ thus generated, along with the $p(t_j)$ from Exercise 1(a), to initialize the Baum-Welch iteration. Compare the resulting plots of average log-probability versus k with those of 1(b), as well as the final values of the average log-probabilities.

Caution: Make sure you mitigate numerical underflow problems when computing the forward and backward probabilities. Use the normalization described in §2.8 if needed.

Submission: Turn in all your plots and discussion, and your source code, via GradeScope; make sure your code is well documented. Points may be deducted for incomprehensible code. Your code may be rerun on different training and test data or with a different initialization to check its correctness; make sure it runs on a linux machine with standard compilers/libraries/environments.