

# Markov Process Lexicon

April 19, 2017

We are interested in modeling the generation of a random lexicon as a Markov process. For the process, we imagine stringing characters together until, at some random point, the string is considered a word. There are thus four qualitatively different states for the Markov process (and degeneracies corresponding to the number of allowable characters). We will need to check

1. the columns sum to 1 (so the transition to the next state is a probability),
2. the largest eigenvalue is equal to 1,
3. the lexicon generated is finite (finite words are generated).

The possible states are the top node (before adding the first letter), adding a character that does not finish a word, adding a character that finishes a word, and the bottom node (corresponding to a finished word that has no suffixes). For  $N$  characters, the transition matrix  $P$  for this process has a shape of  $2(N+1) \times 2(N+1)$  (2 edge nodes,  $N$  characters that don't finish words,  $N$  characters that could finish words). If we imagine each character as equally likely and completable words being completed with probability  $\alpha$ , then the transition matrix can be expressed as

$$P_{ij} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 1/N & 1/2N & \cdots & (1-\alpha)/2N & \cdots & 0 \\ \vdots & \vdots & & \vdots & & \vdots \\ 1/N & 1/2N & \cdots & (1-\alpha)/2N & \cdots & 0 \\ 0 & 1/2N & \cdots & (1-\alpha)/2N & \cdots & 0 \\ \vdots & \vdots & & \vdots & & \vdots \\ 0 & 1/2N & \cdots & (1-\alpha)/2N & \cdots & 0 \\ 0 & 0 & \cdots & \alpha & \cdots & 1 \end{bmatrix}. \quad (1)$$

We've got to unpack this a little. The first row represents the top node state - no state should be able to return to the top node, so it is filled with zeros. The next  $N$  rows correspond to the  $N$  characters that don't end words. The top node leads only to these to

exclude words that are only one character long, but it can lead to any of the characters with equal probability, hence the  $1/N$  populating that segment of the first column. The next  $N$  rows are the  $N$  characters that end words. Characters that do and do not end words may transition to each other (e.g.,  $\text{oend}$  might transition to  $\text{oendi}$  on its way to  $\text{oending\bar{s}}$ ). The latter selects the next character with equal probability ( $1/2N$ ), but the former must retain the possibility of actually ending the branch (as no word can be constructed from  $\text{oending\bar{s}\bullet}$ ). We denote the probability of ending a branch as  $\alpha$ . This  $\alpha$  factor is removed uniformly from the transition probabilities away from a character that end words, hence  $(1 - \alpha)/2N$ . Finally, the last row corresponds to the bottom node state, that cannot transition to anywhere.

The columns can similarly be partitioned into those four categories, where we define the following vectors

$$\mathbf{v_T} = \begin{bmatrix} 0 \\ 1/N \\ \vdots \\ 1/N \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \quad \mathbf{v_{NW}} = \begin{bmatrix} 0 \\ 1/2N \\ \vdots \\ 1/2N \\ 1/2N \\ \vdots \\ 1/2N \\ 0 \end{bmatrix} \quad \mathbf{v_W} = \begin{bmatrix} 0 \\ (1 - \alpha)/2N \\ \vdots \\ (1 - \alpha)/2N \\ (1 - \alpha)/2N \\ \vdots \\ (1 - \alpha)/2N \\ \alpha \end{bmatrix} \quad \mathbf{v_B} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \quad (2)$$

Clearly there are many linear dependencies between the columns (as there are  $N$  each of  $\mathbf{v_{NW}}$  and  $\mathbf{v_W}$ ). An additional linear dependence that is of note is

$$\mathbf{v_W} = \alpha \mathbf{v_B} + (1 - \alpha) \mathbf{v_{NW}} \quad (3)$$

We proceed now to evaluate how each of these vectors is acted upon by  $P$ .

$$P_{ij} \mathbf{v_T} = \begin{bmatrix} 0 \\ N(1/N)(1/2N) \\ \vdots \\ N(1/N)(1/2N) \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1/2N \\ \vdots \\ 1/2N \\ 0 \end{bmatrix} = \mathbf{v_{NW}} \quad (4)$$

$$P_{ij} \mathbf{v_{NW}} = \begin{bmatrix} 0 \\ N(1/2N)^2 + N(1/2N)((1 - \alpha)/2N) \\ \vdots \\ N(1/2N)^2 + N(1/2N)((1 - \alpha)/2N) \\ N(1/2N)\alpha \end{bmatrix} = \begin{bmatrix} 0 \\ (2 - \alpha)/4N \\ \vdots \\ (2 - \alpha)/4N \\ \alpha/2 \end{bmatrix} \quad (5)$$

$$= \frac{2 - \alpha}{2} \mathbf{v_{NW}} + \frac{\alpha}{2} \mathbf{v_B} \quad (6)$$

$$P_{ij}\mathbf{v}_W = \begin{bmatrix} 0 \\ N(1/2N)((1-\alpha)/2N) + N((1-\alpha)/2N)^2 \\ \vdots \\ N(1/2N)((1-\alpha)/2N) + N((1-\alpha)/2N)^2 \\ N\alpha(1-\alpha)/2N + \alpha \end{bmatrix} = \begin{bmatrix} 0 \\ (1-\alpha)(2-\alpha)/4N \\ \vdots \\ (1-\alpha)(2-\alpha)/4N \\ \alpha(3-\alpha)/2 \end{bmatrix} \quad (7)$$

$$= \frac{2-\alpha}{2}\mathbf{v}_W + \frac{\alpha}{2}\mathbf{v}_B \quad (8)$$

$$P_{ij}\mathbf{v}_B = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} = \mathbf{v}_B \quad (9)$$

Note that  $\mathbf{v}_T$  transitions into the “not word” state, which is what we wanted. Next, observe that the bottom node state transitions only to itself, and is therefore our first eigenvector, and has an eigenvalue of one.

$$\mathbf{v}_B \cdot \mathbf{v}_B = 1 \quad (10)$$

$$\mathbf{v}_T \cdot \mathbf{v}_T = 1/N \quad (11)$$

$$\mathbf{v}_{NW} \cdot \mathbf{v}_{NW} = 1/2N \quad (12)$$

$$\mathbf{v}_W \cdot \mathbf{v}_W = (1 - 2\alpha + (1 + 2N)\alpha^2)/2N \quad (13)$$

$$\mathbf{v}_B \cdot \mathbf{v}_T = 0 \quad (14)$$

$$\mathbf{v}_B \cdot \mathbf{v}_{NW} = 0 \quad (15)$$

$$\mathbf{v}_B \cdot \mathbf{v}_W = \alpha \quad (16)$$

$$\mathbf{v}_T \cdot \mathbf{v}_{NW} = 1/2N \quad (17)$$

$$\mathbf{v}_T \cdot \mathbf{v}_W = (1 - \alpha)/2N \quad (18)$$

$$\mathbf{v}_{NW} \cdot \mathbf{v}_W = (1 - \alpha)/2N \quad (19)$$

$$(20)$$

$$\mathbf{u}_1 = \mathbf{v}_B \quad (21)$$

$$\mathbf{u}_2 = \mathbf{v}_W - \frac{\mathbf{v}_B \cdot \mathbf{v}_W}{\mathbf{v}_B \cdot \mathbf{v}_B} \mathbf{v}_B = \mathbf{v}_W - \alpha \mathbf{v}_B \quad (22)$$

$$P_{ij}\mathbf{u}_2 = P(\mathbf{v}_W - \alpha \mathbf{v}_B) = P\mathbf{v}_W - \alpha P\mathbf{v}_B \quad (23)$$

Substituting from eqs 43 and 44

$$= \frac{2-\alpha}{2} \mathbf{v}_W + \frac{\alpha}{2} \mathbf{v}_B - \alpha \mathbf{v}_B \quad (24)$$

$$= \frac{2-\alpha}{2} \mathbf{v}_W - \frac{\alpha}{2} \mathbf{v}_B \quad (25)$$

Try

$$\mathbf{v}_2 = \mathbf{v}_B - \mathbf{v}_{NW} = \begin{bmatrix} 0 \\ -1/2N \\ \vdots \\ -1/2N \\ -1/2N \\ \vdots \\ -1/2N \\ 1 \end{bmatrix} \quad (26)$$

$$P_{ij} \mathbf{v}_2 = \begin{bmatrix} 0 \\ N(-1/2N)(1/2N) + N(-1/2N)((1-\alpha)/2N) \\ \vdots \\ N(-1/2N)(1/2N) + N(-1/2N)((1-\alpha)/2N) \\ N(-1/2N)\alpha + 1 \end{bmatrix} \quad (27)$$

$$= \begin{bmatrix} 0 \\ ((2-\alpha)/2)(-1/2N) \\ \vdots \\ ((2-\alpha)/2)(-1/2N) \\ (2-\alpha)/2 \end{bmatrix} = \frac{2-\alpha}{2} \begin{bmatrix} 0 \\ -1/2N \\ \vdots \\ -1/2N \\ 1 \end{bmatrix} = \frac{2-\alpha}{2} \mathbf{v}_2 \quad (28)$$

$$\mathbf{v}_2 \cdot \mathbf{v}_2 = 2N \frac{1}{4N^2} + 1 = \frac{1+2N}{2N} \quad (29)$$

Alternatively,

$$P \mathbf{v}_2 = P(\mathbf{v}_B - \mathbf{v}_{NW}) = \mathbf{v}_B - \left( \frac{2-\alpha}{2} \mathbf{v}_{NW} + \frac{\alpha}{2} \mathbf{v}_B \right) \quad (30)$$

$$= \frac{2-\alpha}{2} (\mathbf{v}_B - \mathbf{v}_{NW}) = \frac{2-\alpha}{2} \mathbf{v}_2 \quad (31)$$

Now we may write our starting state  $\mathbf{v}_{NW}$  as

$$\mathbf{v}_{NW} = \mathbf{v}_B - \mathbf{v}_2 \quad (32)$$

Now, since

$$P\mathbf{v}_{\mathbf{NW}} = \mathbf{v}_{\mathbf{B}} - \frac{2-\alpha}{2}\mathbf{v}_{\mathbf{2}}, \quad (33)$$

repeated application of the transition matrix yields

$$P^n\mathbf{v}_{\mathbf{NW}} = \mathbf{v}_{\mathbf{B}} - \left(\frac{2-\alpha}{2}\right)^n \mathbf{v}_{\mathbf{2}}. \quad (34)$$

As long as  $0 < \alpha \leq 1$ , we know that  $(2-\alpha)/2 < 1$ , implying that

$$\lim_{n \rightarrow \infty} P^n\mathbf{v}_{\mathbf{NW}} = \mathbf{v}_{\mathbf{B}}, \quad (35)$$

which means that all our branches must eventually end and so we may expect finite words.

## A Errata

We are interested in modeling the generation of a random lexicon as a Markov process. For the process, we imagine stringing characters together until, at some random point, the string is considered a word. There are thus four qualitatively different states for the Markov process (and degeneracies corresponding to the number of allowable characters). We will need to check

1. the columns sum to 1 (so the transition to the next state is a probability),
2. the largest eigenvalue is equal to 1,
3. the lexicon generated is finite (finite words are generated).

The possible states are the top node (before adding the first letter), adding a character that does not finish a word, adding a character that finishes a word, and the bottom node (corresponding to a finished word that has no suffixes). For  $N$  characters, the transition matrix  $P$  for this process has a shape of  $2(N+1) \times 2(N+1)$  (2 edge nodes,  $N$  characters that don't finish words,  $N$  characters that could finish words). If we imagine each character as equally likely and completable words being completed with probability  $\alpha$ , then the transition matrix can be expressed as

$$P_{ij} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 1/N & 1/2N & \cdots & (1-\alpha)/2N & \cdots & 0 \\ \vdots & \vdots & & \vdots & & \vdots \\ 1/N & 1/2N & \cdots & (1-\alpha)/2N & \cdots & 0 \\ 0 & 1/2N & \cdots & (1-\alpha)/2N & \cdots & 0 \\ \vdots & \vdots & & \vdots & & \vdots \\ 0 & 1/2N & \cdots & (1-\alpha)/2N & \cdots & 0 \\ 0 & 0 & \cdots & \alpha & \cdots & 1 \end{bmatrix}. \quad (36)$$

We've got to unpack this a little. The first row represents the top node state - no state should be able to return to the top node, so it is filled with zeros. The next  $N$  rows correspond to the  $N$  characters that don't end words. The top node leads only to these to exclude words that are only one character long, but it can lead to any of the characters with equal probability, hence the  $1/N$  populating that segment of the first column. The next  $N$  rows are the  $N$  characters that end words. Characters that do and do not end words may transition to each other (e.g.,  $\text{oend}\bar{\text{i}}$  might transition to  $\text{oend}\bar{\text{i}}$  on its way to  $\text{oending}\bar{\text{s}}$ ). The latter selects the next character with equal probability ( $1/2N$ ), but the former must retain the possibility of actually ending the branch (as no word can be constructed from  $\text{oending}\bar{\text{s}}\bullet$ ). We denote the probability of ending a branch as  $\alpha$ . This  $\alpha$  factor is removed uniformly from the transition probabilities away from a character that end words, hence  $(1 - \alpha)/2N$ . Finally, the last row corresponds to the bottom node state, that cannot transition to anywhere.

The columns can similarly be partitioned into those four categories, where we define the following vectors

$$\mathbf{v_T} = \begin{bmatrix} 0 \\ 1/N \\ \vdots \\ 1/N \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \quad \mathbf{v_{NW}} = \begin{bmatrix} 0 \\ 1/2N \\ \vdots \\ 1/2N \\ 1/2N \\ \vdots \\ 1/2N \\ 0 \end{bmatrix} \quad \mathbf{v_W} = \begin{bmatrix} 0 \\ (1 - \alpha)/2N \\ \vdots \\ (1 - \alpha)/2N \\ (1 - \alpha)/2N \\ \vdots \\ (1 - \alpha)/2N \\ \alpha \end{bmatrix} \quad \mathbf{v_B} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \quad (37)$$

Clearly there are many linear dependencies between the columns (as there are  $N$  each of  $\mathbf{v_{NW}}$  and  $\mathbf{v_W}$ ). An additional linear dependence that is of note is

$$\mathbf{v_W} = \alpha \mathbf{v_B} + (1 - \alpha) \mathbf{v_{NW}} \quad (38)$$

We proceed now to evaluate how each of these vectors is acted upon by  $P$ .

$$P_{ij} \mathbf{v_T} = \begin{bmatrix} 0 \\ N(1/N)(1/2N) \\ \vdots \\ N(1/N)(1/2N) \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1/2N \\ \vdots \\ 1/2N \\ 0 \end{bmatrix} = \mathbf{v_{NW}} \quad (39)$$

$$P_{ij}\mathbf{v}_{\mathbf{NW}} = \begin{bmatrix} 0 \\ N(1/2N)^2 + N(1/2N)((1-\alpha)/2N) \\ \vdots \\ N(1/2N)^2 + N(1/2N)((1-\alpha)/2N) \\ N(1/2N)\alpha \end{bmatrix} = \begin{bmatrix} 0 \\ (2-\alpha)/4N \\ \vdots \\ (2-\alpha)/4N \\ \alpha/2 \end{bmatrix} \quad (40)$$

$$= \frac{2-\alpha}{2}\mathbf{v}_{\mathbf{NW}} + \frac{\alpha}{2}\mathbf{v}_{\mathbf{B}} \quad (41)$$

$$P_{ij}\mathbf{v}_{\mathbf{W}} = \begin{bmatrix} 0 \\ N(1/2N)((1-\alpha)/2N) + N((1-\alpha)/2N)^2 \\ \vdots \\ N(1/2N)((1-\alpha)/2N) + N((1-\alpha)/2N)^2 \\ N\alpha(1-\alpha)/2N + \alpha \end{bmatrix} = \begin{bmatrix} 0 \\ (1-\alpha)(2-\alpha)/4N \\ \vdots \\ (1-\alpha)(2-\alpha)/4N \\ \alpha(3-\alpha)/2 \end{bmatrix} \quad (42)$$

$$= \frac{2-\alpha}{2}\mathbf{v}_{\mathbf{W}} + \frac{\alpha}{2}\mathbf{v}_{\mathbf{B}} \quad (43)$$

$$P_{ij}\mathbf{v}_{\mathbf{B}} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} = \mathbf{v}_{\mathbf{B}} \quad (44)$$

Note that  $\mathbf{v}_{\mathbf{T}}$  transitions into the “not word” state, which is what we wanted. Next, observe that the bottom node state transitions only to itself, and is therefore our first eigenvector, and has an eigenvalue of one.

$$\mathbf{v}_{\mathbf{B}} \cdot \mathbf{v}_{\mathbf{B}} = 1 \quad (45)$$

$$\mathbf{v}_{\mathbf{T}} \cdot \mathbf{v}_{\mathbf{T}} = 1/N \quad (46)$$

$$\mathbf{v}_{\mathbf{NW}} \cdot \mathbf{v}_{\mathbf{NW}} = 1/2N \quad (47)$$

$$\mathbf{v}_{\mathbf{W}} \cdot \mathbf{v}_{\mathbf{W}} = (1 - 2\alpha + (1 + 2N)\alpha^2)/2N \quad (48)$$

$$\mathbf{v}_{\mathbf{B}} \cdot \mathbf{v}_{\mathbf{T}} = 0 \quad (49)$$

$$\mathbf{v}_{\mathbf{B}} \cdot \mathbf{v}_{\mathbf{NW}} = 0 \quad (50)$$

$$\mathbf{v}_{\mathbf{B}} \cdot \mathbf{v}_{\mathbf{W}} = \alpha \quad (51)$$

$$\mathbf{v}_{\mathbf{T}} \cdot \mathbf{v}_{\mathbf{NW}} = 1/2N \quad (52)$$

$$\mathbf{v}_{\mathbf{T}} \cdot \mathbf{v}_{\mathbf{W}} = (1 - \alpha)/2N \quad (53)$$

$$\mathbf{v}_{\mathbf{NW}} \cdot \mathbf{v}_{\mathbf{W}} = (1 - \alpha)/2N \quad (54)$$

$$(55)$$

$$\mathbf{u}_1 = \mathbf{v}_B \quad (56)$$

$$\mathbf{u}_2 = \mathbf{v}_W - \frac{\mathbf{v}_B \cdot \mathbf{v}_W}{\mathbf{v}_B \cdot \mathbf{v}_B} \mathbf{v}_B = \mathbf{v}_W - \alpha \mathbf{v}_B \quad (57)$$

$$P_{ij} \mathbf{u}_2 = P(\mathbf{v}_W - \alpha \mathbf{v}_B) = P\mathbf{v}_W - \alpha P\mathbf{v}_B \quad (58)$$

Substituting from eqs 43 and 44

$$= \frac{2-\alpha}{2} \mathbf{v}_W + \frac{\alpha}{2} \mathbf{v}_B - \alpha \mathbf{v}_B \quad (59)$$

$$= \frac{2-\alpha}{2} \mathbf{v}_W - \frac{\alpha}{2} \mathbf{v}_B \quad (60)$$

Try

$$\mathbf{v}_2 = \mathbf{v}_B - \mathbf{v}_{NW} = \begin{bmatrix} 0 \\ -1/2N \\ \vdots \\ -1/2N \\ -1/2N \\ \vdots \\ -1/2N \\ 1 \end{bmatrix} \quad (61)$$

$$P_{ij} \mathbf{v}_2 = \begin{bmatrix} 0 \\ N(-1/2N)(1/2N) + N(-1/2N)((1-\alpha)/2N) \\ \vdots \\ N(-1/2N)(1/2N) + N(-1/2N)((1-\alpha)/2N) \\ N(-1/2N)\alpha + 1 \end{bmatrix} \quad (62)$$

$$= \begin{bmatrix} 0 \\ ((2-\alpha)/2)(-1/2N) \\ \vdots \\ ((2-\alpha)/2)(-1/2N) \\ (2-\alpha)/2 \end{bmatrix} = \frac{2-\alpha}{2} \begin{bmatrix} 0 \\ -1/2N \\ \vdots \\ -1/2N \\ 1 \end{bmatrix} = \frac{2-\alpha}{2} \mathbf{v}_2 \quad (63)$$

$$\mathbf{v}_2 \cdot \mathbf{v}_2 = 2N \frac{1}{4N^2} + 1 = \frac{1+2N}{2N} \quad (64)$$



Alternatively,

$$P\mathbf{v}_2 = P(\mathbf{v}_B - \mathbf{v}_{NW}) = \mathbf{v}_B - \left(\frac{2-\alpha}{2}\mathbf{v}_{NW} + \frac{\alpha}{2}\mathbf{v}_B\right) \quad (65)$$

$$= \frac{2-\alpha}{2}(\mathbf{v}_B - \mathbf{v}_{NW}) = \frac{2-\alpha}{2}\mathbf{v}_2 \quad (66)$$