

Consider a random sequence of the 26 characters of the standard Latin alphabet. We want to know the probability $p(n)$ of the sequence AB appearing after the first n characters have been written.

For $n = 0$, $p(0) = 0$.

For $n = 1$, $p(1) = 0$ because we need at least two characters to have been written.

For $n = 2$, $p(2) = 1/26^2$ because of the $1/26$ chance of the first character being A and the $1/26$ chance of the second character being B .

For $n = 3$, $p(3) = 2/26^2$ because either the first two characters are AB , which has a $p(2)$ chance, or the second two characters are AB , which has a $1/26^2$ chance again.

For $n = 4$, $p(4) = 3/26^2 - 1/26^4$. This is because we have the $p(3)$ chance of AB being somewhere in the first 3 characters, or a $1/26^2$ chance of the last two characters being AB . But then we have to subtract the overlap: there is a $1/26^4$ chance that AB is somewhere in the first 3 characters *and* AB is the last two characters (specifically, the sequence $ABAB$ has a $1/26^4$ chance).

For $n = 5$, we can use the same reasoning: there is a $p(4)$ chance of AB being somewhere in the first 4 characters, and there is a $1/26^2$ chance of AB being the fourth and fifth characters respectively. So far that gives us $p(4) + 1/26^2$. But then we have to subtract the overlap again: specifically, the probability that the final two characters are AB and the that AB appears somewhere in the first four characters. If that were the case, then AB would have to actually appear somewhere in the first three characters, since there is no way the fourth character could be B if the sequence ends with AB . Therefore, we have to subtract the overlap $(1/26^2)p(3)$. So $p(5) = p(4) + 1/26^2 - (1/26^2)p(3)$.

In general, this reasoning gives us the recursive sequence

$$p(n) = p(n - 1) + (1/676)(1 - p(n - 2)).$$

We can use linear algebra to examine this sequence.

$$\underbrace{\begin{pmatrix} p(n) \\ p(n - 1) \end{pmatrix}}_{\vec{p}(n)} = \underbrace{\begin{pmatrix} 1 & -1/676 \\ 1 & 0 \end{pmatrix}}_A \underbrace{\begin{pmatrix} p(n - 1) \\ p(n - 2) \end{pmatrix}}_{\vec{p}(n-1)} + \underbrace{\begin{pmatrix} 1/676 \\ 0 \end{pmatrix}}_{\vec{v}},$$

where $\vec{p}(1) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$. Then

$$\vec{p}(2) = A\vec{p}(1) + \vec{v} = \vec{v}$$

$$\vec{p}(3) = A^2\vec{p}(1) + A\vec{v} + \vec{v} = A\vec{v} + \vec{v}$$

$$\vec{p}(4) = A^3\vec{p}(1) + A^2\vec{v} + A\vec{v} + \vec{v} = A^2\vec{v} + A\vec{v} + \vec{v},$$

etc. So really

$$\vec{p}(n) = A^{n-1}\vec{v} + \dots + A\vec{v} + \vec{v} = \sum_{i=0}^{n-1} A^i\vec{v}.$$

This is a geometric sum:

$$\sum_{i=0}^{n-1} A^i = \frac{A^n - I}{A - I} = (A^n - I)(A - I)^{-1},$$

so

$$\vec{p}(n) = (A^n - I)(A - I)^{-1}\vec{v}.$$

We can use this formula to find $\lim_{t \rightarrow \infty} \vec{p}(t)$. Since both of the eigenvalues of A are in the interval $(0, 1)$ we have $\lim_{t \rightarrow \infty} A^t = 0$. So

$$\lim_{t \rightarrow \infty} \vec{p}(t) = (-I)(A - I)^{-1}\vec{v} = \begin{pmatrix} 676 & -1 \\ 676 & 0 \end{pmatrix} \begin{pmatrix} 1/676 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

And so $\lim_{t \rightarrow \infty} p(t) = 1$.

Let's compare this to the probability $q(n)$ of the sequence AA appearing in the first n characters. Then $q(0) = q(1) = 0$ as before, and $q(2) = 1/676$. To calculate $q(3)$, we can take $q(2)$ and add the probability that the last two characters are AA , which is $1/676$. Then we need to subtract the overlap: the event that all three characters are A , which is $1/26^3$. So $q(3) = 2/26^2 - 1/26^3$.

Now let's look at $q(4)$. I think we can take $q(3)$ and add the $1/26^2$ chance that the final two characters are both A , then subtract the overlap which I think is $q(3)(1/26^2)$, but that doesn't feel entirely right because I'm not sure those are independent events. But if so, we get the recursive pattern

$$q(n) = q(n-1) + 1/676 - q(n-1)(1/676) = q(n-1)(675/676) + 1/676,$$

starting with $q(1) = 0$, which is obviously not right (fails for $n = 3$).

In order to get on the right track, it will be instrumental to calculate the probability $Q(n)$ of AA appearing somewhere in the sequence of n characters AND the final character is A , in terms of q . That is,

$$Q(n) = q(n-1)/26 + 1/26^2 - Q(n-1)/26.$$

And

$$q(n) = q(n-1) + 1/26^2 - Q(n-1)/26.$$

So

$$-Q(n-1)/26 = q(n) - q(n-1) - 1/26^2,$$

and so

$$Q(n) = q(n-1)/26 + q(n) - q(n-1),$$

so

$$Q(n-1) = q(n-1) - 25q(n-2)/26,$$

and so finally

$$q(n) = q(n-1) + 1/26^2 - q(n-1)/26 + 25q(n-2)/26^2$$

or

$$q(n) = 25q(n-1)/26 + 25q(n-2)/26^2 + 1/26^2.$$

Now we can set up

$$\underbrace{\begin{pmatrix} q(n) \\ q(n-1) \end{pmatrix}}_{\vec{q}(n)} = \underbrace{\begin{pmatrix} 25/26 & 25/26^2 \\ 1 & 0 \end{pmatrix}}_B \underbrace{\begin{pmatrix} q(n-1) \\ q(n-2) \end{pmatrix}}_{\vec{q}(n-1)} + \underbrace{\begin{pmatrix} 1/26^2 \\ 0 \end{pmatrix}}_{\vec{u}},$$

where $\vec{q}(1) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$. Thus we have the same geometric sum as before:

$$\vec{q}(n) = (B^n - I)(B - I)^{-1}\vec{u}.$$

Again, both of the eigenvalues of B are within $(0, 1)$, so $\lim_{t \rightarrow \infty} B^t = 0$, so

$$\lim_{t \rightarrow \infty} \vec{q}(t) = (-I)(B - I)^{-1}\vec{u} = \begin{pmatrix} 676 & 25 \\ 676 & 26 \end{pmatrix} \begin{pmatrix} 1/676 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix},$$

so $\lim_{t \rightarrow \infty} q(t) = 1$. Look at this graph.