Stochastic processes and Markov chains (part II)

Wessel van Wieringen w.n.van.wieringen@vu.nl

Department of Epidemiology and Biostatistics, VUmc & Department of Mathematics, VU University Amsterdam, The Netherlands



VU medisch centrum



DNA copy number of a genomic segment is simply the number of copies of that segment present in the cell under study.

Healthy normal cell: chr 1 : 2

```
chr 22 : 2
chr X : 1 or 2
chr Y : 0 or 1
```



Chromosomes of a tumor cell



Technique: SKY

The DNA copy number is often categorized into:

- L : *loss* : < 2 copies
- N : normal : 2 copies
- G : *gain* : > 2 copies

In cancer:

- The number of DNA copy number aberrations accumulates with the progression of the disease.
- DNA copy number aberrations are believed to be irreversible.

Let us model the accumulation process of DNA copy number aberrations.

State diagram for the accumulation process of a locus.



The associated initial distribution:

 $\pi = (0, 1, 0)^T$

and, associated transition matrix:

$$\mathbf{P} = \begin{pmatrix} 1 & 0 & 0 \\ \alpha & 1 - \alpha - \beta & \beta \\ 0 & 0 & 1 \end{pmatrix}$$

with parameter constraints:

 $0<\alpha+\beta<1,\alpha>0,\beta>0$

Calculate the probability of a loss, normal and gain at this locus after *p* generations:

$$\begin{pmatrix} P(X_p = L) \\ P(X_p = N) \\ P(X_p = G) \end{pmatrix} = \boldsymbol{\pi}^T \mathbf{P}^p = \begin{pmatrix} (1 - \alpha - \beta)^p \\ \alpha \sum_{t=0}^{p-1} (1 - \alpha - \beta)^t \\ \beta \sum_{t=0}^{p-1} (1 - \alpha - \beta)^t \end{pmatrix}$$

Using:

$$\sum_{t=0}^{\infty} (1-c)^t = \frac{1}{c}$$

These probabilities simplify to, e.g.:

$$P(X_p = L) = \frac{\alpha}{\alpha + \beta} [1 - (1 - \alpha - \beta)^p]$$

In practice, a sample is only observed once the cancer has already developed. Hence, the number of generations *p* is unknown. This may be accommodated by modeling *p* as being Poisson distributed:

$$P(Y = p) = \lambda^p \exp(-\lambda)/p!$$

This yields, e.g.:

$$P(X = N) = \sum_{p=0}^{\infty} P(X = L | Y = p) P(Y = p)$$
$$= \sum_{p=0}^{\infty} (1 - \alpha - \beta)^p \lambda^p \exp(-\lambda)/p!$$

So far, we only considered one locus. Hence:



For multiple loci:



Multiple loci \rightarrow multivariate problem.

Complications:

- *p* unknown,
- loci not independent.

Solution:

- p random,
- assume particular dependence structure.

After likelihood formulation and parameter estimation:

- identify most aberrated loci,
- reconstruct time of onset of cancer.

Stationary distribution

We generated a DNA sequence of 10.000.000 bases long in accordance with a 1st order Markov chain.

For stretches DNA ever longer and ever farther away from the first base pair we calculated the nucleotide %.

bp	bp	% A	%C	%G	8 T	
001	5e+03	0.014	0.069	0.262	0.655	
001	2e+04	0.045	0.239	0.228	0.488	1
001	1e+05	0.144	0.319	0.211	0.327	C C
001	5e+05	0.158	0.295	0.205	0.342	en
001	2e+06	0.142	0.284	0.216	0.357	bre
•	•	•	•	•	•	nve
•	•	•	•	•	•	Ö
•	•	•	•	•	•	¥
		0.150	0.280	0.220	0.350	

Stationary distribution

Hence, after a while an "equilibrium" sets in. Not necessarily a fixed state or pattern, but:

the proportion of a time period that is spent in a particular state converges to a limit value.

The limit values of all states form the *stationary distribution* of the Markov process, denoted by:

$$\boldsymbol{\varphi} = (\varphi_1, \varphi_2, \dots, \varphi_S)^T \quad \text{with } \sum_{k=1}^S \varphi_k = 1$$

For a stationary process, it holds that $P(X_t = E_i) = \varphi_i$ for all *t* and *i*.

In particular:

$$P(X_t = E_i) = \phi_i = P(X_{t+1} = E_i).$$

Of course, this does *not* imply:

$$\mathsf{P}(X_{t}=E_{i}, X_{t+1}=E_{i}) = \varphi_{i} \varphi_{i}$$

as this ignores the 1st order Markov dependency.

Stationary distribution

The stationary distribution is associated with the firstorder Markov process, parameterized by (π, P) .

Question How do ϕ and (π, P) relate?

Hereto, recall:

$$\varphi_i = P(X_t = E_i)$$

$$\varphi_i = P(X_{t+1} = E_i)$$

definition

$$\varphi_i = P(X_{t+1} = E_i)$$

$$= \sum_{k=1}^{S} P(X_{t+1} = E_i, X_t = E_k)$$

use the fact that $P(A, B) + P(A, B^{C}) = P(A)$

$$\varphi_{i} = P(X_{t+1} = E_{i})$$

$$= \sum_{k=1}^{S} P(X_{t+1} = E_{i}, X_{t} = E_{k})$$

$$= \sum_{k=1}^{S} P(X_{t+1} = E_{i} | X_{t} = E_{k}) P(X_{t} = E_{k})$$

use the definition of conditional probability: P(A, B) = P(A | B) P(B)

$$\varphi_i = P(X_{t+1} = E_i)$$

$$= \sum_{k=1}^{S} P(X_{t+1} = E_i, X_t = E_k)$$

$$= \sum_{k=1}^{S} P(X_{t+1} = E_i | X_t = E_k) P(X_t = E_k)$$

$$= \sum_{k=1}^{S} p_{ki} \varphi_k$$

$$\varphi_i = P(X_{t+1} = E_i)$$

$$= \sum_{k=1}^{S} P(X_{t+1} = E_i, X_t = E_k)$$

$$= \sum_{k=1}^{S} P(X_{t+1} = E_i | X_t = E_k) P(X_t = E_k)$$

$$= \sum_{k=1}^{S} p_{ki} \varphi_k$$

Thus:

$$\varphi^T = \varphi^T \mathbf{P} \longrightarrow \mathsf{Eigenvectors!}$$

Stationary distribution

Theorem

Irreducible, aperiodic Markov chains with a finite state space *S* have a stationary distribution to which the chain converges as $t \rightarrow \infty$.



A Markov chain is *aperiodic* if there is no state that can only be reached at multiples of a certain period. E.g., state E_i only at t = 0, 3, 6, 9, et cetera.

Example of an aperiodic Markov chain





...ATCGATCGATCGATCG...

Example of a periodic Markov chain



The fraction of time spent in **A** (roughly ϕ_A):

$$P(X_{t+1000} = A) = \frac{1}{4}$$

whereas
$$P(X_{t+1000} = \mathbf{A} | X_4 = \mathbf{A}) = 1$$
.

Stationary distribution

A Markov chain is *irreducible* if every state can (in principle) be reached (after enough time has passed) from every other state.

Examples of a reducible Markov chain



C is an absorbing state.



A will not be reached.

How do we find the stationary distribution?

We know the stationary distribution satisfies:

$$arphi^T = arphi^T \mathbf{P}$$
 (*)
and
 $arphi_1 + arphi_2 + \ldots + arphi_S = 1$

We thus have S+1 equations for S unknowns: one of the equations in (*) can be dropped (which is irrelevant), and the system of S remaining equations needs to be solved.

Consider the transition matrix:

$$\mathbf{P} = \begin{pmatrix} 0.35 & 0.65 \\ 0.81 & 0.19 \end{pmatrix}$$

In order to find the stationary distribution we need to solve the following system of equations:

$$\begin{array}{rcl} 0.35\,\varphi_1 + 0.81\,\varphi_2 &=& \varphi_1 \\ & \varphi_1 + \varphi_2 &=& 1 \end{array}$$

This yields: $(\phi_1, \phi_2)^T = (0.5547945, 0.4452055)^T$

On the other hand, for *n* large:

$$P(X_{t+n} = E_j \mid X_t = E_i) = \varphi_j$$

is independent of *i*. Or, $p_{ij}^{(n)} = \phi_j$.

Hence, the *n*-step transition matrix $\mathbf{P}^{(n)}$ has identical rows:

$$\lim_{n \to \infty} \mathbf{P}^{(n)} = \lim_{n \to \infty} \mathbf{P}^n = \begin{pmatrix} \varphi_1 & \varphi_2 & \dots & \varphi_S \\ \varphi_1 & \varphi_2 & \dots & \varphi_S \\ \vdots & \vdots & & \vdots \\ \varphi_1 & \varphi_2 & \dots & \varphi_S \end{pmatrix}$$

This motivates a numerical way to find the stationary distribution.

Same example as before:

$$\mathbf{P} = \begin{pmatrix} 0.35 & 0.65 \\ 0.81 & 0.19 \end{pmatrix}$$

with stationary distribution: $(\phi_1, \phi_2)^T = (0.5547945, 0.4452055)^T$

Then:

$$\mathbf{P}^{(2)} = \begin{pmatrix} 0.35 & 0.65 \\ 0.81 & 0.19 \end{pmatrix} \begin{pmatrix} 0.35 & 0.65 \\ 0.81 & 0.19 \end{pmatrix}$$
$$= \begin{pmatrix} 0.35 \times 0.35 + 0.65 \times 0.81 & 0.35 \times 0.65 + 0.65 \times 0.19 \\ 0.81 \times 0.35 + 0.19 \times 0.81 & 0.81 \times 0.65 + 0.19 \times 0.19 \end{pmatrix}$$

matrix multiplication ("rows times columns")

Thus:

$$\mathbf{P}^{(2)} = \begin{pmatrix} 0.6490 & 0.3510 \\ 0.4374 & 0.5626 \end{pmatrix}$$

In similar fashion we obtain:

$$\mathbf{P}^{(5)} = \begin{pmatrix} 0.5456249 & 0.4543751 \\ 0.5662212 & 0.4337788 \end{pmatrix}$$

$$\mathbf{P}^{(20)} = \begin{pmatrix} 0.5547946 & 0.4452054 \\ 0.5547944 & 0.4452056 \end{pmatrix}$$

Convergence of the stationary distribution

Convergence to the stationary distribution

Define the vector $\mathbf{1} = (1, \dots, 1)^T$.

We have already seen:

 $\mathbf{P}^{n} = \mathbf{1} \, \boldsymbol{\varphi}^{\mathsf{T}} \qquad \text{for large } n$

Question

How fast does \mathbf{P}^n go to $\mathbf{1} \boldsymbol{\varphi}^T$ as $n \to \infty$?

Answer

- 1) Use linear algebra
- 2) Calculate numerically

Convergence to the stationary distribution

Fact

The transition matrix **P** of a finite, aperiodic, irreducible Markov chain has an eigenvalue equal to 1 (λ_1 =1), while all other eigenvalues are (in the absolute sense) smaller than one: $|\lambda_k| < 1$, k=2,...,3.

Focus on λ_1 =1

We know $\mathbf{\phi}^{\mathsf{T}} \mathbf{P} = \mathbf{\phi}^{\mathsf{T}}$ for the stationary distribution. Hence, $\mathbf{\phi}$ is the left eigenvector of eigenvalue $\lambda = 1$.

Also, row sums of **P** equal **1**: **P 1** = **1**. Hence, **1** is a right eigenvector of eigenvalue λ =1.

The spectral decomposition of a square matrix P is given by: $\mathbf{P} = \mathbf{V} \mathbf{D} \mathbf{V}^{-1}$

where:

- D diagonal matrix containing the eigenvalues,
- $\mathbf V$ columns contain the corresponding eigenvectors.

In case of **P** is symmetric, **V** is orthogonal: $\mathbf{V}^{-1} = \mathbf{V}^T$ Then:

$$\mathbf{P} = \mathbf{V} \mathbf{D} \mathbf{V}^T$$

Convergence to the stationary distribution

The *spectral decomposition* of **P**, reformulated:



The eigenvectors are normalized:

$$\mathbf{r}_k^T \boldsymbol{\ell}_k = 1, \quad \text{and} \\ \mathbf{r}_{k_1}^T \boldsymbol{\ell}_{k_2} = 0 \quad \text{if } k_1 \neq k_2$$


We now obtain the spectral decomposition of the *n*-step transition matrix \mathbf{P}^n . Hereto observe that:

$$\mathbf{P}^{n} \mathbf{r}_{k} = \mathbf{P}^{n-1} \Big(\sum_{k_{0}=1}^{S} \lambda_{k_{0}} \mathbf{r}_{k_{0}} \boldsymbol{\ell}_{k_{0}}^{T} \Big) \mathbf{r}_{k}$$

plug in the spectral decomposition of ${\bf P}$

We now obtain the spectral decomposition of the *n*-step transition matrix \mathbf{P}^n . Hereto observe that:

$$\mathbf{P}^{n} \mathbf{r}_{k} = \mathbf{P}^{n-1} \Big(\sum_{k_{0}=1}^{S} \lambda_{k_{0}} \mathbf{r}_{k_{0}} \boldsymbol{\ell}_{k_{0}}^{T} \Big) \mathbf{r}_{k}$$

$$= \mathbf{P}^{n-1} \Big(\sum_{k_{0}=1}^{S} \lambda_{k_{0}} \mathbf{r}_{k_{0}} [\mathbf{r}_{k}^{T} \boldsymbol{\ell}_{k_{0}}]^{T} \Big)$$

bring the right eigenvector in the sum, and use the properties of the transpose operator

C

We now obtain the spectral decomposition of the *n*-step transition matrix \mathbf{P}^n . Hereto observe that:

$$\mathbf{P}^{n} \mathbf{r}_{k} = \mathbf{P}^{n-1} \Big(\sum_{k_{0}=1}^{S} \lambda_{k_{0}} \mathbf{r}_{k_{0}} \boldsymbol{\ell}_{k_{0}}^{T} \Big) \mathbf{r}_{k}$$

$$= \mathbf{P}^{n-1} \Big(\sum_{k_{0}=1}^{S} \lambda_{k_{0}} \mathbf{r}_{k_{0}} [\mathbf{r}_{k}^{T} \boldsymbol{\ell}_{k_{0}}]^{T} \Big)$$

$$= \mathbf{P}^{n-1} \lambda_{k} \mathbf{r}_{k}$$

the eigenvectors are normalized

C

We now obtain the spectral decomposition of the *n*-step transition matrix \mathbf{P}^n . Hereto observe that:

$$\mathbf{P}^{n} \mathbf{r}_{k} = \mathbf{P}^{n-1} \Big(\sum_{k_{0}=1}^{S} \lambda_{k_{0}} \mathbf{r}_{k_{0}} \boldsymbol{\ell}_{k_{0}}^{T} \Big) \mathbf{r}_{k}$$

$$= \mathbf{P}^{n-1} \Big(\sum_{k_{0}=1}^{S} \lambda_{k_{0}} \mathbf{r}_{k_{0}} [\mathbf{r}_{k}^{T} \boldsymbol{\ell}_{k_{0}}]^{T} \Big)$$

$$= \mathbf{P}^{n-1} \lambda_{k} \mathbf{r}_{k}$$

$$= \lambda_{k} \mathbf{P}^{n-1} \mathbf{r}_{k}$$

Repeating this argument *n* times yields:

$$\mathbf{P}^{n} \mathbf{r}_{k} = \lambda_{k}^{n} \mathbf{r}_{k}$$
$$\boldsymbol{\ell}_{k}^{T} \mathbf{P}^{n} = \lambda_{k}^{n} \boldsymbol{\ell}_{k}^{T}$$

Hence, ℓ_k and \mathbf{r}_k are left and right eigenvector with eigenvalue λ_k^n of \mathbf{P}^n . Thus:

$$\mathbf{P}^n = \sum_{k=1}^S \lambda_k^n \mathbf{r}_k \boldsymbol{\ell}_k^T$$



Verify the spectral decomposition for \mathbf{P}^2 :

$$\mathbf{P}^{2} = \left(\sum_{k=1}^{S} \lambda_{k} \mathbf{r}_{k} \boldsymbol{\ell}_{k}^{T}\right) \left(\sum_{k=1}^{S} \lambda_{k} \mathbf{r}_{k} \boldsymbol{\ell}_{k}^{T}\right)$$
$$= \sum_{k_{1}=1}^{S} \sum_{k_{2}=1}^{S} \lambda_{k_{1}} \mathbf{r}_{k_{1}} \boldsymbol{\ell}_{k_{1}}^{T} \lambda_{k_{2}} \mathbf{r}_{k_{2}} \boldsymbol{\ell}_{k_{2}}^{T}$$
$$= \sum_{k_{1}=1}^{S} \sum_{k_{2}=1}^{S} \lambda_{k_{1}} \lambda_{k_{2}} \mathbf{r}_{k_{1}} \boldsymbol{\ell}_{k_{1}}^{T} \mathbf{r}_{k_{2}} \boldsymbol{\ell}_{k_{2}}^{T}$$
$$= \sum_{k_{1}=1}^{S} \sum_{k_{2}=1}^{S} \lambda_{k_{1}} \lambda_{k_{2}} \mathbf{r}_{k_{1}} (\mathbf{r}_{k_{2}}^{T} \boldsymbol{\ell}_{k_{1}})^{T} \boldsymbol{\ell}_{k_{2}}^{T}$$
$$= \sum_{k_{1}=1}^{S} \lambda_{k_{1}}^{2} \mathbf{r}_{k} \boldsymbol{\ell}_{k}^{T}$$

Use the spectral decomposition of \mathbf{P}^n to show how fast \mathbf{P}^n converges to $\mathbf{1} \ \boldsymbol{\phi}^T$ as $n \to \infty$.

We know:

 λ_1 =1, $|\lambda_k|$ < 1 for k=2, ...,S, $\ell_1 = arphi$ and $\mathbf{r}_1 = \mathbf{1}$

Then:

$$\mathbf{P}^{n} = 1^{n} \mathbf{1} \boldsymbol{\varphi}^{T} + \sum_{k=2}^{S} \lambda_{k}^{n} \mathbf{r}_{k} \boldsymbol{\ell}_{k}^{T}$$

Expanding this: $\mathbf{P}^{n} = \mathbf{1}^{n} \mathbf{1} \boldsymbol{\varphi}^{T} + \sum_{k=2}^{S} \lambda_{k}^{n} \mathbf{r}_{k} \boldsymbol{\ell}_{k}^{T}$ $= \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} (\varphi_1, \varphi_2, \dots, \varphi_S) + \lambda_2^n \mathbf{r}_2 \boldsymbol{\ell}_2^T + \dots + \lambda_S^n \mathbf{r}_S \boldsymbol{\ell}_S^T$

Clearly:

$$\lim_{n \to \infty} \mathbf{P}^n = \begin{pmatrix} \varphi_1 & \varphi_2 & \dots & \varphi_S \\ \varphi_1 & \varphi_2 & \dots & \varphi_S \\ \vdots & \vdots & & \vdots \\ \varphi_1 & \varphi_2 & \dots & \varphi_S \end{pmatrix}$$

Furthermore, as:

 $0 \le |\lambda_{k_1}| < |\lambda_{k_2}| < 1 \qquad \Rightarrow \qquad 0 \le |\lambda_{k_1}^n| < |\lambda_{k_2}^n| < 1$

It is the second largest (in absolute sense) eigenvalue that dominates, and thus determines the convergence speed to the stationary distribution.

Fact

A Markov chain with a symmetric **P** has a uniform stationary distribution.

Proof

- Symmetry of **P** implies that left- and right eigenvectors are identical (up to a constant).
- First right eigenvector corresponds vector of ones, **1**.
- Hence, the left eigenvector equals c1.
- The left eigenvector is the stationary distribution and should sum to one: c = 1 / (number of states).

Question

Suppose the DNA may reasonably described by a first order Markov model with transition matrix **P**:

$$\mathbf{P} = \begin{pmatrix} 0.2 & 0.3 & 0.3 & 0.2 \\ 0.1 & 0.4 & 0.4 & 0.1 \\ 0.3 & 0.2 & 0.2 & 0.3 \\ 0.4 & 0.1 & 0.1 & 0.4 \end{pmatrix}$$

and stationary distribution:

$$\varphi^T = (\varphi_A, \varphi_C, \varphi_G, \varphi_T) = (1/4, 1/4, 1/4, 1/4)$$

and eigenvalues:

$$\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \lambda_3, \lambda_4)^T = (1, 0.2, 0, 0)^T$$

Question

What is the probability of a **G** at position 2, 3, 4, 5, 10? And how does this depend on the starting nucleotide?

In other words, give:

$$P(X_{2} = G | X_{1} = A) = \dots$$

$$P(X_{2} = G | X_{1} = C) = \dots$$

$$P(X_{2} = G | X_{1} = G) = \dots$$

$$P(X_{2} = G | X_{1} = T) = \dots$$

But also:

$$P(X_3 = G | X_1 = A) = ...$$

et cetera.

Thus, calculate $P(X_t = G | X_1 = x_1)$ with *t*=2, 3, 4, 5, 10, and $x_1 = A$, C, G, T.

x1=A x1=C x1=G x1=T t= 2 0.300000 0.400000 0.200000 0.100000 t= 3 0.2600000 0.280000 0.2400000 0.2200000 t= 4 0.2520000 0.2560000 0.2480000 0.2440000 t= 5 0.2504000 0.2512000 0.2496000 0.2488000 t=10 0.2500001 0.2500004 0.2499999 0.2499996

Study the influence of the first nucleotide on the calculated probability for increasing *t*.



```
> # define \pi and P
> pi <- matrix(c(1, 0, 0, 0), ncol=1)
> P <- matrix(c(2, 3, 3, 2, 1, 4, 4, 1, 3, 2, 2, 3,
     4, 1, 1, 4), ncol=4, byrow=TRUE)/10
> # define function that calculates the powers of a
> # matrix (inefficiently though)
> matrixPower <- function(X, power) {</pre>
     Xpower <-X
     for (i in 2:power) {
           Xpower <- Xpower %*% X
     return (Xpower)
  }
```

```
> # calculate P to the power 100
> matrixPower(P, 100)
```

Question

Suppose the DNA may reasonably described by a first order Markov model with transition matrix **P**:

Р	=	(0.77450	0.22500	0.00025	0.00025
			0.22500	0.77450	0.00025	0.00025
			0.00025	0.00025	0.77450	0.22500
			0.00025	0.00025	0.22500	0.77450

and stationary distribution:

$$\varphi^T = (\varphi_A, \varphi_C, \varphi_G, \varphi_T) = (1/4, 1/4, 1/4, 1/4)$$

and eigenvalues:

 $\boldsymbol{\lambda} = (1, 0.9990, 0.5495, 0.5495)^T$

Again calculate $P(X_t = G | X_1 = x_1)$ with *t*=2, 3, 4, 5, 10, and $x_1 = A$, C, G, T.



Now the influence of the first nucleotide fades slowly. This can be explained by the large 2nd eigenvalue. Processes back in time So far, we have studied Markov chains forward in time. In practice, we may wish to study processes back in time.

Example

Evolutionary models that describe occurrence of SNPs in DNA sequences. We aim to attribute two DNA sequences to a common ancestor.



Consider a Markov chain $\{X_t\}_{t=1,2,...}$. The reverse Markov chain $\{X_r^*\}_{r=1,2,...}$ is then defined by:

$$X_{r}^{*} = X_{t-r}$$



With transition probabilities:

$$p_{ij} = P(X_t = E_j | X_{t-1} = E_i)$$

 $p_{ij}^* = P(X_r^* = E_j | X_{r-1}^* = E_i)$

$$p_{ij}^* = P(X_r^* = E_j | X_{r-1}^* = E_i)$$

just the definition

$$p_{ij}^{*} = P(X_{r}^{*} = E_{j} | X_{r-1}^{*} = E_{i})$$

$$= P(X_{t-r} = E_{j} | X_{t-r+1} = E_{i})$$

express this in terms of the original Markov chain using that $X_r^* = X_{t-r}$

$$p_{ij}^{*} = P(X_{r}^{*} = E_{j} | X_{r-1}^{*} = E_{i})$$

$$= P(X_{t-r} = E_{j} | X_{t-r+1} = E_{i})$$

$$= P(X_{t-r+1} = E_{i} | X_{t-r} = E_{j}) \frac{P(X_{t-r} = E_{j})}{P(X_{t-r+1} = E_{i})}$$

apply definition of conditional probability twice (Bayes): P(A | B) = P(B | A) P(A) / P(B)

$$p_{ij}^{*} = P(X_{r}^{*} = E_{j} | X_{r-1}^{*} = E_{i})$$

= $P(X_{t-r} = E_{j} | X_{t-r+1} = E_{i})$
= $P(X_{t-r+1} = E_{i} | X_{t-r} = E_{j}) \frac{P(X_{t-r} = E_{j})}{P(X_{t-r+1} = E_{i})}$
= $p_{ji} \varphi_{j} / \varphi_{i}$

Hence:

$$p_{ij}^* = p_{ji} \frac{\varphi_j}{\varphi_i}$$

Check that rows of the transition matrix **P*** sum to one, i.e.:

$$p_{i1}^{*} + p_{i2}^{*} + \dots + p_{jS}^{*} = 1$$

Hereto:

The two Markov chains defined by **P** and **P*** have the same stationary distribution. Indeed, as:

$$\sum_{k=1}^{S} \varphi_k p_{kj}^* = \sum_{k=1}^{S} p_{ik} \varphi_j$$
$$= \varphi_j \sum_{k=1}^{S} p_{ik}$$
$$= \varphi_j$$

we have:

$$arphi^T = arphi^T \mathbf{P}^*$$

Processes back in time

Definition

A Markov chain is called **reversible** if $p_{ij}^* = p_{ij}$. In that case:

$$p_{ij}^* = p_{ij} = p_{ji} \varphi_j / \varphi_i$$

Or,

$$\varphi_i p_{ij} = \varphi_j p_{ji}$$
 for all i and j.

These are the so-called *detailed balance equations*.

Theorem

A Markov chain is reversible if and only if the detailed balance equations hold.

Example 1

The 1st order Markov chain with transition matrix:

$$\mathbf{P} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$$

is irreversible. Check that this (deterministic) Markov chain does not satisfy the detailed balance equations.

Irreversibility can be seen from a sample of this chain: ... ABCABCABCABCABC...

In the reverse direction transitions from **B** to **C** do not occur!

Example 2

The 1st order Markov chain with transition matrix:

$$\mathbf{P} = \begin{pmatrix} 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \\ 0.8 & 0.1 & 0.1 \end{pmatrix}$$

is irreversible. Again, a uniform stationary distribution: $\varphi^T = (\varphi_A, \varphi_B, \varphi_C) = (1/3, 1/3, 1/3)$

As **P** is not symmetric, the detailed balance equations are not satisfies:

$$p_{ij} / 3 \neq p_{ji} / 3$$
 for all i and j.

Example 2 (continued)

The irreversibility of this chain implies:

$$P(\mathbf{A} \rightarrow \mathbf{B} \rightarrow \mathbf{C} \rightarrow \mathbf{A})$$

$$= P(\mathbf{A} \rightarrow \mathbf{B}) P(\mathbf{B} \rightarrow \mathbf{C}) P(\mathbf{C} \rightarrow \mathbf{A})$$

$$= 0.8 * 0.8 * 0.8$$

$$\neq 0.1 * 0.1 * 0.1$$

$$= P(\mathbf{A} \rightarrow \mathbf{C}) P(\mathbf{C} \rightarrow \mathbf{B}) P(\mathbf{B} \rightarrow \mathbf{A})$$

$$= P(\mathbf{A} \rightarrow \mathbf{C} \rightarrow \mathbf{B} \rightarrow \mathbf{A}).$$

It matters how one walks from **A** to **A**.

Or, it matters whether one walks forward or backward.

Kolmogorov condition for reversibility

A stationary Markov chain is reversible if and only if any path from state E_i to state E_i has the same probability as the path in the opposite direction. Or, the Markov chain is reversible if and only if:

 $p_{i,i_1} \cdot p_{i_1,i_2} \cdot \ldots \cdot p_{i_k,i} = p_{i,i_k} \cdot \ldots \cdot p_{i_2,i_1} \cdot p_{i_1,i_2}$ for all *i*, *i*₁, *i*₂, ..., *i*_k.

E.g.:
$$P(A \rightarrow B \rightarrow C \rightarrow A) = P(A \rightarrow C \rightarrow B \rightarrow A).$$

Interpretation

For a reversible Markov chain it is not possible to determine the direction of the process from the observed state sequence alone.

- Molecular phylogenetics aims to reconstruct evolutionary relationships between present day species from their DNA sequences. Reversibility is then an essential assumption.
- Genes are transcribed in one direction only (from the 3' end to the 5' end). The promoter is only on the 3' end. This suggests irreversibility.

E. Coli

For a gene in the E. Coli genome, we estimate:

```
Transition matrix
```

[,1] [,2] [,3] [,4] [1,] 0.2296984 0.3155452 0.2273782 0.2273782 [2,] 0.1929134 0.2421260 0.2933071 0.2716535 [3,] 0.1979167 0.2854167 0.2166667 0.3000000 [4,] 0.2522686 0.2032668 0.2341198 0.3103448

Stationary distribution [1] 0.2187817 0.2578680 0.2436548 0.2796954

Then, the detailed balance equation do not hold, e.g.: $\pi_1 p_{12} \neq \pi_2 p_{21.}$

Note

Within evolution theory the notion of irreversibility refers to the presumption that complex organisms once lost evolution will not appear in the same form.

Indeed, the likelihood of reconstructing a particular phylogenic system is infinitesimal small.

Application: motifs

Study the sequence of the promotor region up-stream of a gene.



This region contains binding sites for the transcription factors that regulate the transcription of the gene.



The binding sites of a transciption factor (that may regulate multiple genes) share certain sequence patterns, *motifs*.

Not all transcription factors and motifs are known. Hence, a high occurrence of a particular sequence pattern in the upstream regions of a gene may indicate that it has a regulatory function (e.g., binding site).

Problem

Determine the probability of observing *m* motifs in a background generated by a 1st order stationary Markov chain.
An *h*-letter **word** $W = w_1 w_2 \dots w_h$ is a map from $\{1, \dots, h\}$ to \mathcal{A}^h , where \mathcal{A} some non-empty set, called the **alphabet**.

In the DNA example:

$$\mathcal{A} = \{ \texttt{A, C, G, T} \}$$

and, e.g.:

- W = CAGTACGACT
- W = TACGACTGCATATGCGTA

A word *W* is *p*-periodic if

 $w_{t_1} = w_{t_2}$ for all $t_1 \equiv t_2 \mod p$, $1 \le t_1, t_2 \le p$

The lag between two overlapping occurrences of the word.

The set of all periods of W (less than h) is the *period set*, denoted by $\mathcal{P}(W)$. In other words, the set of integers 0 < p< h such that a new occurrence of W can start p letters after an occurrence of W.

```
If W_1 = CGATCGATCG, then \mathcal{P}(W) = \{5, 9\}
For:
```

```
123456789
CGATCGATC
CGATCGATC
CGATCGATC
```

If W_2 = CGAACTG, then $\mathcal{P}(W) = \emptyset$

Let N(W) be the number of (overlapping) occurrences of an *h*-letter word W in a random sequence n on A.

If Y_t is the random variable defined by:

 $Y_t = I_{\{\text{an occurrence of } W \text{ starts at position } t\}}$

then

$$N(W) = \sum_{t=1}^{n-h+1} Y_t$$

Also, denote the number of occurrences in W by:

$$n(a) = \{\#t \mid a \in W\}$$

and

$$n(a\bullet) = \sum_{b \in \mathcal{A}} n(ab)$$

If W = CATGA, then: n(A) = 2 and n(A•) = 1 Assume a stationary 1^{st} order Markov model for the random sequence of length *n*. The probability of an occurrence of *W* in the random sequence is given by:

$$\mu(W) = \varphi_{w_1} \prod_{t=1}^{h-1} (\mathbf{P})_{w_j, w_{j+1}}$$

In the 1^{st} order Markov model, the expected number of occurrences of W is approximated by:

$$\hat{N}(W) = \frac{N(w_1w_2) \cdot N(w_2w_3) \cdot \ldots \cdot N(w_{h-1}w_h)}{N(w_2) \cdot N(w_3) \cdot \ldots \cdot N(w_{h-1})}$$

and its variance by:

$$\sigma^{2}(W) = \mu(W) + \sum_{p \in \mathcal{P}(W)} \mu(w_{1} \dots w_{p} w_{1} \dots w_{h}) + \mu(W)^{2} \Big(\sum_{a \in \mathcal{A}} \frac{n(a \bullet)^{2}}{\mu(a)} - \sum_{a_{1}, a_{2} \in \mathcal{A}} \frac{n(a_{1}a_{2})^{2}}{\mu(a_{1}a_{2})} + \frac{1 - n(w_{1} \bullet)^{2}}{\mu(w_{1})} \Big)$$

To find words with exceptionally frequency in the DNA, the following (asymptotically) standard normal statistic is used:

$$Z(W) = \frac{N(W) - \hat{N}(W)}{\sqrt{n\,\hat{\sigma}^2(W)}} \sim N(0,1)$$

The *p*-value of word *W* is then given by:

$$p(W) = P(Z \ge Z(W))$$

= $1 - \Phi_{0,1}(Z(W))$

Note

Robin, Daudon (1999) provide exact probabilities of word occurences in random sequences.

However, Robin, Schbath (2001) point out that calculation of the exact probabilities is computationally intensive. Hence, the use of an approximation here. References & further reading

References and further reading

- Ewens, W.J, Grant, G (2006), *Statistical Methods for Bioinformatics*, Springer, New York.
- Reinert, G., Schbath, S., Waterman, M.S. (2000), "Probabilistic and statistical properties of words: an overview", *Journal of Computational Biology*, **7**, 1-46.
- Robin, S., Daudin, J.-J. (1999), "Exact distribution of word occurrences in a random sequence of letters", *Journal of Applied Probability*, **36**, 179-193.
- Robin, S., Schbath, S. (2001), "Numerical comparison of several approximations of the word count distribution in random sequences", *Journal of Computational Biology*, **8**(4), 349-359.
- Schbath, S. (2000), "An overview on the distribution of word counts in Markov chains", *Journal of Computational. Biology*, **7**, 193-202.
- Schbath, S., Robin, R. (2009), "How can pattern statistics be useful for DNA motif discovery?". In *Scan Statistics: Methods and Applications* by Glaz, J. *et al.* (eds.).



This material is provided under the Creative Commons Attribution/Share-Alike/Non-Commercial License.

See http://www.creativecommons.org for details.