Stochastic processes and Markov chains (part I)

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



Example 1

- The intensity of the sun.
- Measured every day by the KNMI.
- Stochastic variable X_t represents the sun's intensity at day t, 0 ≤ t ≤ T. Hence, X_t assumes values in R⁺ (positive values only).



Example 2

- DNA sequence of 11 bases long.
- At each base position there is an **A**, **C**, **G** or **T**.
- Stochastic variable X_i is the base at position i, i = 1,...,11.
- In case the sequence has been observed, say:

$$(x_1, x_2, \ldots, x_{11}) = \text{ACCCGATAGCT},$$

then **A** is the realization of X_1 , **C** that of X_2 , et cetera.



Stochastic processes



. . .

. . .

Example 3

- A patient's heart pulse during surgery.
- Measured continuously during interval [0, 7].
- Stochastic variable X_t represents the occurrence of a heartbeat at time t, $0 \le t \le T$. Hence, X_t assumes only the values 0 (no heartbeat) and 1 (heartbeat).



Figure: http://www.gnkaterini.gr/

Example 4

- Brain activity of a human under experimental conditions.
- Measured continuously during interval [0, *T*].
- Stochastic variable X_t represents the magnetic field at time t, 0 ≤ t ≤ T. Hence, X_t assumes values on R.



The state space S is the collection of values that the random variables of the stochastic process may assume.

If $S = \{E_1, E_2, ..., E_s\}$, then X_t is a discrete stochastic variable.

If $S = [0, \infty)$, then X_t is a continuous stochastic variable.

Time can also be either discrete or continuous.



Ewens, Grant (2005):

Section 4.5.

book

6

Stochastic process:

- Discrete time: t = 0, 1, 2, 3,
- $S = \{-3, -2, -1, 0, 1, 2, \ldots\}$
- P(one step up) = $\frac{1}{2}$ = P(one step down)
- The process is locked in at -3



The *first passage time* of a certain state E_i in *S* is the time t at which $X_t = E_i$ for the first time since the start of the process.

The *absorbing state* is a state E_i for which the following holds: if $X_t = E_i$ than $X_s = E_i$ for all $s \ge t$. The process will never leave state E_i .

The *time of absorption* of an absorbing state is the first passage time of that state.



A stochastic process is described by a collection of time points, the state space and the simultaneous distribution of the variables X_t , i.e., the distributions of all X_t and their dependency.

There are two important types of processes:

- *Poisson process*: all variables are identically and independently distributed. Examples: queues for counters, call centers, servers, et cetera.
- *Markov process*: the variables are dependent in a simple manner.



Markov processes

book Axelson-Fisk (2010): Definition 2.1.

A 1st order Markov process in discrete time is a stochastic process $\{X_t\}_{t=1,2,...}$ for which the following holds:

$$P(X_{t+1}=x_{t+1} \mid X_t=x_t, \dots, X_1=x_1) = P(X_{t+1}=x_{t+1} \mid X_t=x_t).$$

In other words, only the present determines the future, the past is irrelevant.

From every state with probability 0.5 one step up or down (except for state -3).



The 1st order Markov property, i.e.

$$P(X_{t+1}=x_{t+1} | X_t=x_t, ..., X_1=x_1) = P(X_{t+1}=x_{t+1} | X_t=x_t),$$

does not imply independence between X_{t-1} and X_{t+1} .

Nor does it imply that $P(X_{t+1}=x_{t+1} | X_{t-1}=x_{t-1})$ equals zero.



The Oth order Markov property, i.e.

$$P(X_{t+1}=x_{t+1} \mid X_t=x_t, \ldots, X_1=x_1) = P(X_{t+1}=x_{t+1}).$$

Thus, X_t and X_{t+1} are independent. So are X_{t-1} and X_{t+1} , and X_t and X_{t+100} , and ...

Or:

$$P(X_{t+1}=x_{t+1}, X_t=x_t) = P(X_{t+1}=x_{t+1}) * P(X_t=x_t)$$

Classic example:

- Tossing a coin : $P(X_{t+1} = head | X_t = tail) = P(X_{t+1} = head)$
- Rolling dice : $P(X_{t+1} = 6 | X_t = 6) = P(X_{t+1} = 6)$

Difference between 0^{th} and 1^{st} order Markov process \rightarrow DNA example

- \rightarrow state space: {**A**, **C**, **G**, **T**}
- \rightarrow sample in accordance with a 0th and 1st Markov chain
- Recall: 1st order takes into account the previous state visited, 0th order does not.



Markov processes

Question: what is the order of the processes?

ATTE TCTACCACTCTATAAAATCAA T C TACC TTCC AT AT AC CACC CACCTC ATA TTCCC CT ATC CCAC A C TCC TTA C C CCAA CTTCTT AC TTACTAC FCATTET T ATTEATATECCE CE A CE A CAACATAC T CTT TET T TECA TETACTET TETACACAA CEA A CEATE TETETAATE C ATA A AT T T TT CACTT ACC ATATAT AA A TC CCTCA CTC ATCTTA TCTAT CCAT T TTCCTTTCA ATACATCCATA C AACA CACCACTTA AACCCA TA CTAT C CCAT ICAATCO TO AATACTCCTATACTTACATTACAAACCTITO TAT CA A ACA AC AC T AT CTCCC AC CCC TO TAT A AATT AACCTTC TC TCCAA C AACC ACA CTCACA TA AA TC AATTCTT AT C C AA ACACCTTCA ACTCA TTTT CACCTTAA C CC AACO T TITT T C CTCTTTTAT TTATCAAC A TCC CCT TO TITTO ACCTARCE CCCCA. ATAT T ARC TECCETTAR AT TATCE ATAR. A AAT AT ATTECT TA TE AAAC CCAACTAC TECAA TA C TETECECA TEATEC CT TA A A CCACA CACTA TCC TT C C TA C CAC AACTAAACT A TTACTCA TTTT TACTTT T AC TTA TCCTCTACAATAA AACTT T T TC TATACAC TA A A T TICC ATTC T ACCCITCICITC TTACITCA CCCTCT TCTTTCC CAATTCTTTCCACATTAAAA A AAATA AA CC ATACA AA ACC TCATACCACAAC CAC TAT FTCC TTT CTCTATA C AA TACHA AATCC C T CCTAT C CCATACCTAAC AT CCCTA T CTC TCA AA ATATTAT CTTTTCA AA CAA A CT TAAACTAACA ATO C T T C C CCAAA T TCTCT CAATCA CC TT ACTITATTA ACC AATAATTCC C TAA CT C AACCTC C TCA AAC C TCACAACTTACAAA CTACACTCCCACAT TCAAATCTAACATCT TTC CTTC CT CALLA AA CCTCAAACT ACAAC TCATTTACT CAT CAATTTA TCCTT TA ACC AC T CAAAC TAA CC TCTC CCAACTT CCAAA AT ATA AACTA TTT CCCAC CC CCT CATTAA C C A TT ATCACAT ATCCTA TITCTICAC TC AAT ACC AC CCC TICT TCACTAA T T T CCCT/ T AC TTC TTA CC TC TCTC T TCC CA TAT T ATTATT TCACATCATTTCC CTA C TC CCT AT AAAA TC C T CCC TTAACAC T CTAC AACTTC C T C FACCATCAATC CCCA AACC ATCCATTT T TAC CAACATC CCACCCAT TAT / CC CTITECT C ATCC CAACCA C TCT ACCT TCATA TAC CT CTA ACTT TTCT CCTCCATTT CACT T T CCCC CTTT CAAATTCACA A A T TT ACT TC CA AC AACC A TC AACT CCA AC AA TCAACTC TCC AACC ATCTATCCCAACCTATTTTC CATC T AC AAA CCATCCTTTA CATTA TCTTAAC TTATA TAAA TA TAAT CCT TT T CCCCCTC AC CACA ACTT TC AACTTAATA CACTAC ATAA CA ACTCTCTACT CCT TA AAA TAA AATA T TT TTC ATTATTA CACTCAC C TO TAATAAC A C ATA TTT A A ACTACAC TTCTAAATC AACCATCTTC ATTTTTT AAACTCATT TAC TAA TTTATCTACCC ATCTTCTTA AT AT CTCTC CTAC A TCCT A TTA CTAAATTT A ATCATCATAAATA CAT TAAT CCAT T C C AATTAT CT ATTCCCC A CCA CTATTATC TITTCCATCC C TAT TTCC AT TCAATTCCTTATT C A TTAATACCT CTC CTT T TCA TITITA / C CC TCTCCAATCA A T CT CAA ACTTCA A CTCCTTCTC T CTATACAC C CALACATAATTA A CATAT T T C T CCCACTCTC CCCCCA CTA TC A CALLT AA TCA T AATCTCCA T ATA AA C TITTCTAAC C CAA T AC ATAT TTA T TCCCCCCAAAACCAAA TTACCCTTA TAA TITC T CACAAC TTAATATTAAC AC CI TT T TCCACC TACATATCAA C CTTCAC A TACACA C CACTC CACAAATACTTCT CCATTIT AA TTA TA CA A CTA TITTACCCT TC CTTCA TCAG TAACCTCC CTCCCCCTCACAAACAC T C CTCT ATC CTTT AAAT TCATCACTTACCATTTAACATCTCAATTC TCTTC ATACAT CACTAC TTCC TATACATA C ACAT ATAAT ICAACTEAATAA TTEECAE A T C C CC TAT ATC TAAATAC A TTAAATA CAACTTT CTT A TA CT UTTEACACCT A A CEEAA T C ACCTAATTEETECEAT TCC TATT A TACTCTCTCTAT CTT TACC TITE E TTAA CE CE CETEACATA ACTT CACTT TE T TEAC TE AA CE T TE E AAACTTA TEAC E T CE TACAA TA ACTCACCO CTCACA CTCA ATTAC ATCTAAC CCAC C TC ACTTAAC CAATCTTACT CCA AAT CT CT CAT AATCA A T TTC CATATA TA ATATTTC TA A AT AACTTCAC ACA AAC T CCTCT CTTCAA ACAT CCAA AAATCC CAATTAATAAAAAA ATAC C TAT CACTACTTCAATA ACA T CCCT CTCC A CTC AATC A T ACA AT AC CTCCCAATCA C CTATACCAT T ATAT CT C CTT TACCCTTATA TAATTC A T CAACC ATAACCAACAAA TC TADAAC CCCATT AA CAT C TATO ACCANATCE TITTEEET A E TA AE ACCAE AAA A AT ETTT ETEE CAE ATA ECTECETE ACAATTTAA T AE T T E AE ATAA ECEE TE TT AAT T ACTAT TC TATCATC ACCOTITCA TATTCCCACCCAA TCTC AACTCTT TT ATTCCC T CC AC A AACCTA ACTAA A A ACTT C CCTCTTTA CCTCT A T TACC AA TCTT TITITE ATT CT. AT TITCTACCTT CTITC T C CTCAAC T. ACAC TT C AAATAT TCCTCCAATCCA. AA TAATC CA C CAAAT CTA C C TAATTA CCTC AT ATTA ICTCT AAATTCCTCT CCT CCTC CCAAC TTC CTAACTACCA CACCA CTACA T TAATTCAC CT A T TA AAT TACT T AAT CCAATCAT AAC T ACC TTCCCCCC TAT T CATCACTA ATAAC AATCACACTT ATTTCCAAC TITT AATTATTACT IT TITTCT AC ATCAA T CC C CTAACA C ATCCA TCAATTTAC CTCAC TCTTT T CTCTTTAAA CC AAC ACCAC A CAAAAT ACA COTITAC AC C TTA ATA TATT CTCC AATAT AT CCC AC AA C TA ACCO CAC C TTAA TC CII TCCAT C CTCCCC TATA CATTTAA C CTCTACCC CCTACCT ACC ACA A C TCC TTTA ACAACTACAA CCTACTCC CAC CC T ATCA TT TACCTTTT A AT ACAT CT T C CCTT AC CCTACT TCCCACATCTCT ACACC TCC C TTATT CA CCC T CAA ACT C ACCCCCAC CC TAC CAACATACACCCAAAAC C CCCAT CTATTTCAC CTC TCC TTATTACTT CAA TTCC CAAATT CTCTTCTCAAACCATC C CCC T TAT CACATA T TA AAA ATCCAAAC TATTACTCCCACTA CATATC CACC AATA AAT CACAT CC T CT C CC AA CATCTTTTCTAATTAAAT TT T TATTTAATACA AC TABLETATASCAT A TOCTO ANTA ATT AACCAACCT TO CA COT TTTTTA TAAC

Markov processes

Question: what is the order of the processes?

T TTTTTTT TT T TTTT DOT DO DO D TT TTTTTTTT CCACCCCACCACCACACAC TTT GROUP TO DO TO TO TO T TT SECTITES TO A STATE TT THE TT TT TOUT SECTOR TO TT TT SECTORS AND A TO T T TITE T. T T TT TOTAL TOTAL TITLE TITLE TO AND A TTAL TOTAL TT TTT T TITITIO TOTITI OF CATGO TOTITIONS IN THIS CONTACT. T T TTTTT T T T TT T T T T T T T T T T TTT TT TTT TT TT TT TTTT T T TITITIT TT TT T TTT T TT TTT AN TITLE THE TRANSPORT TENT TO THE TRANSPORT OF THE TAX TO THE TAX TO THE AT THE AT TAX AND A A A A A A A A A A TTT TT TTTT. T. TTT CTOTAL TOTAL T. T. O. T. T. T. T. T. G. TTTT TTT TT T TTT TTTT OT CTITITIONS TO TITITISTIC TITICOUT CONTROL THE THE TO CTITICE THE OTHER THE TABLES TO THE OFF CTITIES TT T T TTT TTT TT TTTTTTTT T TTO CTTTS TO TO TO TTS TO CTTTS OF TO TT SET TITLE FOR T T TT TITLE TITLE TITLES TO A TABLE TO A TABLE TO A TABLE TO A TABLE TABLE TABLE TO A TABLE T THE T STT TTTTTTTT TT TTTT T OTI Т T TTT CCCACAACCCCCACAACACCCCCCAAAAAAAAAACAACCAACCAACAACAACCACCAACAAACCACC **COTTOTTO** TTTTTGTT TO TTT TTTTG TTGGT TTT T TTT TTT T TTT TTTT T T T TT T TTT T TTT TTTTT TTT OF GRATITATION OF TROCCORD TATITIVET OF TRATES THE TIME T T T TT TT TT TTTTTTT TTT T CONTENT TO THE TOTAL CONTENT OF THIS OF THE TOT TTITTTTT T T TTITTTT T TT THETE OT AN ANTING THIT TT ATTIT TT TTTT TTTTTT TT T TTTT TTTTTTTTTT T COTTO CITITITA TIGA TO C TT T T T T TT TTTTTTTT TT T TT TTO THE TO THE OTHER TO TTTTT TO TTO T TTT T T T T T T TTT T T T T TTT TT T TTT TT ACACCCACCC OGT TO TITTO IT ORDET TITT TITT TT TITT TT TT TT TT OTT OTTITT T TO OTTITT T TO OTT TT TT TT TT OCTOOTICOTOTITTC TLOGT AUGULT T TITUE TIT AU T T TT AUTT AT AUGUTT TIT AUTTIT. T T T TT TT T TT TT TT TT TTT TTT T TTTTT TT TT TT TTO TOTAL TOTAL TTO CTT TT T TTT TTTT T TTTTT T TT TTT TTTT T TT TT TT T TT T T TTT TTTT TTATA TO SO TROOTING THE TO SO THE THE TO SO THEFT THE THE THE TO TATA TO SO TO SO THE THEFT OF A SO THE THEFT TTTTT TT TTTT TTT TTTT TTT T T TTT T TT TTTT T T T T TT TTTT T T THE TABLE TO THE TITLE OF THE TABLE TABLE TO THE TABLE



An *m*-order Markov process in discrete time is a stochastic process $\{X_t\}_{t=1,2,...}$ for which the following holds:

$$P(X_{t+1}=x_{t+1} \mid X_t=x_t, \dots, X_1=x_1) = P(X_{t+1}=x_{t+1} \mid X_t=x_t, \dots, X_{t-m+1}=x_{t-m+1}).$$

Loosely, the future depends on the most recent past.

Suppose *m*=9:



Distribution of the *t*+1-th base depends on the 9 preceeding ones.

A Markov process is called a *Markov chain* if the state space is discrete, i.e., is finite or countable. We consider Markov chains in discrete time.

The *transition probabilities* are the $P(X_{t+1}=x_{t+1} | X_t=x_t)$, but also the

 $P(X_{t+1}=x_{t+1} \mid X_s=x_s)$ for s < t, where x_{t+1}, x_t in $\{E_1, ..., E_s\} = S$.



A Markov process is called *time homogeneous* if the transition probabilities are independent of *t*.

$$P(X_{t+1}=x_1 \mid X_t=x_2) = P(X_{s+1}=x_1 \mid X_s=x_2).$$

For example:

$$P(X_2 = -1 \mid X_1 = 2) = P(X_6 = -1 \mid X_5 = 2)$$

Example of a *time inhomogeneous* process (not in lecture).



Axelson-Fisk (2010):

Definition 2.3.

book

book Ewens, Grant (2005): Section 4.9.

Consider a DNA sequence of 11 bases. Then, $S=\{A, C, G, T\}$, X_t is the base of position t, and $\{X_t\}_{t=1, ..., 11}$ is a Markov chain if the base of position t only depends on the base of position t-1, and not on those before t-1.

If this is plausible, a 1st Markov chain may be an acceptable model for base ordering in DNA sequences.



book Ewens, Grant (2005): Section 4.6.

Denote the transition probabilities of a finite, time homogeneous Markov chain in discrete time $\{X_t\}_{t=1,2,...}$ with $S=\{E_1, ..., E_s\}$ as:

$$P(X_{t+1}=E_j \mid X_t=E_i) = p_{ij} \quad (\text{does not depend on } t).$$

Putting the p_{ii} in a matrix yields the transition matrix:

$$\mathbf{P} = \begin{pmatrix} p_{11} & p_{12} & \cdots & p_{1s} \\ p_{21} & p_{22} & & p_{2s} \\ \vdots & & \ddots & \vdots \\ p_{s1} & p_{s2} & \cdots & p_{ss} \end{pmatrix}$$

The rows of this matrix sum to one.

Coin: $P(X=head) = \frac{1}{2} \rightarrow P(X=tail) = 1 - P(X=head)$

Markov processes



where

$$p_{AA} = P(X_{t+1} = \mathbf{A} \mid X_t = \mathbf{A})$$

and:

$$p_{AA} + p_{AC} + p_{AG} + p_{AT} = 1$$

 $p_{CA} + p_{CC} + p_{CG} + p_{CT} = 1$
et cetera



Markov processes

Question Which state diagram corresponds to the transition matrix **P**?







The *initial distribution* $\mathbf{\pi} = (\pi_1, ..., \pi_s)^T$ gives the probabilities of the initial state:

$$\pi_i = P(X_1 = E_i)$$
 for $i = 1, ..., s$,

and $\pi_1 + \dots + \pi_s = 1$.

Together the initial distribution π and the transition matrix **P** determine the probability distribution of the process, the Markov chain $\{X_t\}_{t=1,2,...}$

Markov processes

Question

Consider the transition matrix **P** (to the right) and initial distribution $\pi = (0, 0.5, 0.5, 0).$

$$\mathbf{P} = \mathbf{C} \begin{pmatrix} \mathbf{A} & \mathbf{C} & \mathbf{G} & \mathbf{T} \\ \mathbf{A} & \mathbf{C} & \mathbf{G} & \mathbf{T} \\ \mathbf{A} & \mathbf{C} & \mathbf{G} & \mathbf{T} \\ 0.1 & 0.9 & 0 & 0 \\ 0 & 0.1 & 0.9 & 0 \\ \mathbf{G} & \mathbf{C} & \mathbf{0} & 0.1 & 0.9 \\ \mathbf{T} & 0.9 & 0 & 0.1 & 0.9 \\ 0.9 & 0 & 0 & 0.1 \end{pmatrix}$$

Calculate the probabilities: $\rightarrow P(X_1=A, X_2=C),$ $\rightarrow P(X_1=C, X_2=G),$ $\rightarrow P(X_1=C, X_2=T).$



We now show how the couple (π, P) determines the probability distribution (transition probabilities) of time steps larger than one.

Ewens, Grant (2005):

Section 4.6.

Hereto define for n = 2: $p_{ij}^{(2)} = P(X_{t+2} = E_j | X_t = E_i)$ general n $p_{ij}^{(n)} = P(X_{t+n} = E_j | X_t = E_i)$

Now express $p_{ii}^{(2)}$ in terms of (π , **P**).

For
$$n = 2$$
:
 $p_{ij}^{(2)} = P(X_{t+2} = E_j | X_t = E_i)$
just the definition

For
$$n = 2$$
:
 $p_{ij}^{(2)} = P(X_{t+2} = E_j | X_t = E_i)$
 $= \sum_{k=1}^{S} P(X_{t+2} = E_j, X_{t+1} = E_k | X_t = E_i)$
 $= \sum_{k=1}^{S} P(X_{t+2} = E_j | X_{t+1} = E_k, X_t = E_i)$
 $\times P(X_{t+1} = E_k | X_t = E_i)$

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

For
$$n = 2$$
:
 $p_{ij}^{(2)} = P(X_{t+2} = E_j | X_t = E_i)$
 $= \sum_{k=1}^{S} P(X_{t+2} = E_j, X_{t+1} = E_k | X_t = E_i)$
 $= \sum_{k=1}^{S} P(X_{t+2} = E_j | X_{t+1} = E_k, X_t = E_i)$
 $\times P(X_{t+1} = E_k | X_t = E_i)$
 $= \sum_{k=1}^{S} P(X_{t+2} = E_j | X_{t+1} = E_k) P(X_{t+1} = E_k | X_t = E_i)$

use the Markov property

For
$$n = 2$$
:
 $p_{ij}^{(2)} = P(X_{t+2} = E_j | X_t = E_i)$
 $= \sum_{k=1}^{S} P(X_{t+2} = E_j, X_{t+1} = E_k | X_t = E_i)$
 $= \sum_{k=1}^{S} P(X_{t+2} = E_j | X_{t+1} = E_k, X_t = E_i)$
 $\times P(X_{t+1} = E_k | X_t = E_i)$
 $= \sum_{k=1}^{S} P(X_{t+2} = E_j | X_{t+1} = E_k) P(X_{t+1} = E_k | X_t = E_i)$
 $= \sum_{k=1}^{S} p_{kj} p_{ik} = \sum_{k=1}^{S} p_{ik} p_{kj} = (\mathbf{P}^2)_{ij}$

book Ewens, Grant (2005): Section 4.6.

In summary, we have shown that:

$$\mathbf{P}^{(2)} = \mathbf{P}^2$$

In a similar fashion we can show that:

$$\mathbf{P}^{(n)} = \mathbf{P}^n \quad \text{for all } n \ge 2.$$

In words: the transition matrix for *n* steps is the onestep transition matrix raised to the power *n*. The general case is proven (see SM) by induction to *n*. This requires the *Kolmogorov-Chapman equations*:

$$(\mathbf{P}^{n+m})_{ij} = \sum_{k=1}^{S} (\mathbf{P}^n)_{ik} (\mathbf{P}^m)_{kj}$$

for all n, m \geq 0 and i, j =1,...,S.

Kolmogorov-Chapman equations illustrated:



A numerical example:

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

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:

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

 $0.6490 = P(X_{t+2} = I | X_t = I)$ is composed of two probabilities:


In similar fashion we may obtain:

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

Sum over probs. of all possible paths between 2 states:



Similarly:

 $\left(\begin{array}{ccc} 0.5547946 & 0.4452054 \\ 0.5547944 & 0.4452056 \end{array}
ight)$ $\mathbf{P}^{(20)} =$

Question: consider the transition matrix:

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

Then $(\mathbf{P}^2)_{1,1}$ corresponds to the probability over paths:



Question: consider the transition matrix:

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

Then $(\mathbf{P}^2)_{1,2}$ corresponds to the probability over paths:



How to sample from a Markov process? Consider the DNA example.

 $\mathbf{A} \quad \mathbf{C} \quad \mathbf{G} \quad \mathbf{T} \\ \boldsymbol{\pi} = (0.45, 0.05, 0.25, 0.25)^T$



$$\mathbf{P} = \mathbf{M} \left\{ \begin{matrix} \mathbf{f} \\ \mathbf{r} \\ \mathbf{O} \\ \mathbf{m} \end{matrix} \right\} \left\{ \begin{matrix} \mathbf{A} \\ \mathbf{C} \\ \mathbf{G} \\ \mathbf{G} \\ \mathbf{T} \end{matrix} \left(\begin{matrix} 0.1 & 0.1 & 0.1 & 0.7 \\ 0.2 & 0.3 & 0.3 & 0.2 \\ 0.4 & 0.3 & 0.2 & 0.1 \\ 0.1 & 0.3 & 0.5 & 0.1 \end{matrix} \right) \right\}$$

Markov process

1 Sample the first base from the initial distribution π : P(X_1 =A) = 0.45, P(X_1 =C) = 0.10, et cetera

Suppose we draw an **A**. Our DNA sequence now consists of one base, namely **A**.

2 Sample the next base using **P** given the previous base. Here $X_1 = \mathbf{A}$, thus we sample from:

 $P(X_2=A | X_1=A) = 0.1, P(X_2=C | X_1=A) = 0.1, et cetera.$

Our DNA sequence after the second step: AT

3 The last step is iterated until the last base. DNA sequence: **ATCCGATGC**

In the famous first application of Markov chains, Andrey Markov studied the sequence of 20,000 letters in Pushkin's poem "Eugeny Onegin", discovering that

- the stationary vowel probability is 0.432,
- the probability of a vowel following a vowel is 0.128, and
- the probability of a vowel following a consonant is 0.663.

Markov also studied the sequence of **100,000** letters in Aksakov's novel "The Childhood of Bagrov, the Grandson".

This was in 1913, long before the computer age!

Estimation (from a sample)



Likelihood

The *likelihood* is a function from the model parameter space, where the parameter (say) θ lives, to the probability space:

 $L: \boldsymbol{\theta} \to [0,1]$

The likelihood yields the probability of the observed data **X** for any parameter choice:

$$L(\mathbf{X};\boldsymbol{\theta}) = P(\mathbf{X}_1,\ldots,\mathbf{X}_n;\boldsymbol{\theta})$$

If the observations are independent, this factorizes to:

$$L(\mathbf{X}; \boldsymbol{\theta}) = \prod_{i=1}^{n} P(\mathbf{X}_i; \boldsymbol{\theta})$$

Likelihood (example)

Let X_i be a random variable representing the outcome of tossing of a coin, either head (H) or tail (T). The typical distributional assumption for X_i is:

 $X \sim \text{Bernoulli}(\theta)$

where $\theta = P(X_i = \mathbf{T})$.

Assume four realizations of X_i are available:

$$\mathbf{x} = (x_1, x_2, x_3, x_3) = (\mathbf{H}, \mathbf{H}, \mathbf{T}, \mathbf{H})$$

Construct the likelihood of for these data.

Likelihood (example) The likelihood is:

 $L(\mathbf{X} = \mathbf{x}; \theta)$ 4 $= P(\mathbf{X} = \mathbf{x}; \theta) = \prod P(X_i = x_i; \theta)$ i=1 $= P(X_1 = \mathbb{H}; \theta) \times P(X_2 = \mathbb{H}; \theta)$ $\times P(X_3 = \mathsf{T}; \theta) \times P(X_4 = \mathsf{H}; \theta)$ $= (1-\theta) \times (1-\theta) \times \theta \times (1-\theta)$ $= \theta (1-\theta)^3.$

Maximum likelihood

The parameter value for which the likelihood attains its maximum is referred to as the *maximum likelihood* estimate:

$$\hat{\boldsymbol{\theta}}_{ML} = \arg \max_{\boldsymbol{\theta}} L(\mathbf{X}; \boldsymbol{\theta})$$

The maximum likelihood estimate is the parameter value that, given the model, is most likely to have given rise to the observed data.

Or loosely, that explains the observations best.

Likelihood (example continued) Return to the 'coin example'. Its likelihood:



Question: How does one find its maximum?

Likelihood (example continued)

Return to the 'coin example'. The maximum likelihood estimate is found by maximizing the likelihood. Due to the concavity of the likelihood, this is equivalent to finding the maximum of the log-likelihood:

$$\log[L(\mathbf{X} = \mathbf{x}; \theta)] = \log(\theta) + 3\log(1 - \theta).$$

Equate the derivative w.r.t. θ to zero:

$$\frac{\partial}{\partial \theta} \log[L(\mathbf{X} = \mathbf{x}; \theta)] = \frac{1}{\theta} - \frac{3}{1 - \theta} = 0.$$

And solve for θ to arrive at:

$$\hat{\boldsymbol{\theta}}_{ML} = \frac{1}{4}.$$

Maximum likelihood estimation

A general procedure, based on an appealing principle, to derive estimators.

ML estimation comprises:

- \rightarrow specification of the distribution of the random variable.
- \rightarrow formulation of the likelihood.
- \rightarrow taking the logarithm of the likelihood (for convenience)
- \rightarrow search for the (location of the) maximum:
 - take derivative with respect to parameters,
 - equate derivative to zero,
 - find zeros of this equation.

Likelihood (example 2)

Let Y_i be a continuous random variable following a normal distribution:

$$Y_i \sim \mathcal{N}(\mu, \sigma^2)$$

Obtain the likelihood for continuous random variables, the density (instead of probability) is used. That of the normal distribution is well-known.

Assuming the independence assumption, the likelihood for a collection of *n* samples is:

$$L(\mathbf{Y} = \mathbf{y}; \mu, \sigma^2) = \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp[-(y_i - \mu)^2/(2\sigma^2)]$$

Likelihood (example 2)

Take the logarithm and obtain the log-likelihood:

$$\sum_{i=1}^{n} \left\{ -\frac{1}{2} \log(2\pi\sigma^2) - (y_i - \mu)^2 / (2\sigma^2) \right\}$$

The log-likelihood is maximized at:



Parameter estimation

Likelihood:

Given the 1st order Markov model, we may write down the likelihood of an observed sequence:

 $P(\texttt{GCTAATTCAGCT}; \mathbf{P})$

The likelihood is a function from the model parameter space, where transition matrix **P** lives, to the probability space [0,1]. It yields the probability of the observed data for any parameter choice.

Note

The likelihood does not factorize as before as contiguous DNA bases are not independent.

The likelihood for the sequence AG:

$$P(X_1 = \mathbf{A}, X_2 = \mathbf{G})$$

= $P(X_2 = \mathbf{G} \mid X_1 = \mathbf{A}) P(X_1 = \mathbf{A})$

The likelihood for the sequence **AGT**:

$$\begin{split} P(X_1 = \mathtt{A}, X_2 = \mathtt{G}, X_3 = \mathtt{T}) \\ &= P(X_3 = \mathtt{T} \mid X_1 = \mathtt{A}, X_2 = \mathtt{G}) \\ &\times P(X_1 = \mathtt{A}, X_2 = \mathtt{G}) \\ &= P(X_3 = \mathtt{T} \mid X_2 = \mathtt{G}) \\ &\times P(X_2 = \mathtt{G} \mid X_1 = \mathtt{A}) \, P(X_1 = \mathtt{A}) \end{split}$$

More general, using the definition of conditional probability, the likelihood can be decomposed as follows:

$$L(\mathbf{X}) = P(X_1 = x_1, X_2 = x_2, \dots, X_T = x_T)$$

= $P(X_T = x_T | X_1 = x_1, \dots, X_{T-1} = x_{T-1})$
 $\times P(X_{T-1} = x_{T-1} | X_1 = x_1, \dots, X_{T-2} = x_{T-2})$

$$\times P(X_2 = x_2 | X_1 = x_1)$$
$$\times P(X_1 = x_1)$$

where $x_1, ..., x_T$ in $S = \{E_1, ..., E_S\}$.

Using the Markov property, we obtain:

$$L(\mathbf{X}) = P(X_1 = x_1, X_2 = x_2, \dots, X_T = x_T)$$

= $P(X_T = x_T | X_{T-1} = x_{T-1})$
 $\times P(X_{T-1} = x_{T-1} | X_{T-2} = x_{T-2})$
 \dots
 $\times P(X_2 = x_2 | X_1 = x_1)$
 $\times P(X_1 = x_1)$
 T

$$= P(X_1 = x_1) \prod_{t=2}^{n} P(X_t = x_t | X_{t-1} = x_{t-1})$$

Question

Consider the 1st order Markov model with initial distribution $\mathbf{\pi} = (0, 0.5, 0.5, 0)$ and transition matrix:



Express the likelihood of the sequence: GCATA in terms of the parameters of the 1st order Markov chain. The indicator function:

$$I_{\{X=\mathbf{A}\}} = \begin{cases} 1 & \text{if } X = \mathbf{A} \\ 0 & \text{if } X \neq \mathbf{A} \end{cases}$$

Recall:

$$\prod_{k=1}^{5} x^{a_k} = x^{a_1} x^{a_2} x^{a_3} x^{a_4} x^{a_5}$$

$$= x^{a_1 + a_2 + a_3 + a_4 + a_5}$$

$$= x^{\sum_{k=1}^{5} a_k}$$

Parameter estimation

Furthermore, e.g.:

$$P(X_{t} = x_{t} | X_{t-1} = x_{t-1})$$

$$= p_{AA}^{I\{X_{t} = A, X_{t-1} = A\}} p_{AC}^{I\{X_{t} = C, X_{t-1} = A\}} p_{AG}^{I\{X_{t} = T, X_{t-1} = A\}} p_{AT}^{I\{X_{t} = T, X_{t-1} = A\}} p_{AT}^{I\{X_{t} = A, X_{t-1} = C\}} p_{CC}^{I\{X_{t} = C, X_{t-1} = C\}} p_{CG}^{I\{X_{t} = T, X_{t-1} = C\}} p_{CT}^{I\{X_{t} = T, X_{t-1} = C\}} p_{CT}^{I\{X_{t} = A, X_{t-1} = C\}} p_{GC}^{I\{X_{t} = C, X_{t-1} = C\}} p_{GG}^{I\{X_{t} = T, X_{t-1} = C\}} p_{GT}^{I\{X_{t} = T, X_{t-1} = C\}} p_{GT}^{I\{X_{t} = A, X_{t-1} = C\}} p_{GC}^{I\{X_{t} = C, X_{t-1} = C\}} p_{GG}^{I\{X_{t} = T, X_{t-1} = C\}} p_{TT}^{I\{X_{t} = A, X_{t-1} = T\}} p_{TC}^{I\{X_{t} = C, X_{t-1} = T\}} p_{TG}^{I\{X_{t} = T, X_{t-1} = T\}} p_{TT}^{I\{X_{t} = T, X_$$

where only one transition probability at the time enters the likelihood, due to the indicator function.

The likelihood then becomes:

$$L(\mathbf{X}) = P(X_1 = x_1) \\ \times p_{AA}^{\sum_{t=2}^{T} I_{\{X_t = A, X_{t-1} = A\}}} \cdot \dots \cdot p_{TT}^{\sum_{t=2}^{T} I_{\{X_t = T, X_{t-1} = T\}}}$$

$$= P(X_1 = x_1) p_{AA}^{n_{AA}} \cdot \ldots \cdot p_{TT}^{n_{TT}}$$

where, e.g.,

$$n_{AA} = \sum_{t=2}^{T} I_{\{X_t = A, X_{t-1} = A\}}$$

Recall that, e.g.,

$$p_{AA} + p_{AC} + p_{AG} + p_{AT} = 1.$$

Or,

$$p_{\text{AT}} = 1 - p_{\text{AA}} - p_{\text{AC}} - p_{\text{AG}}$$

Substitute this in the likelihood, and take the logarithm to arrive at the log-likelihood.

Note: it is irrelevant which transition probability is substituted.

The log-likelihood:

 $\log[L(\mathbf{X})]$ $= \log[P(X_1 = x_1)]$ $+ n_{AA} \log(p_{AA}) + n_{AC} \log(p_{AC})$ $+ n_{AG} \log(p_{AG}) + n_{AT} \log(1 - p_{AA} - p_{AC} - p_{AG})$ \dots $+ n_{TA} \log(p_{TA}) + n_{TC} \log(p_{TC})$

 $+ n_{TG} \log(p_{TG}) + n_{TT} \log(1 - p_{TA} - p_{TC} - p_{TG})$

Differentiation the log-likelihood yields, e.g.:

$$\frac{\partial \log[L(\mathbf{X})]}{\partial p_{AA}} = \frac{n_{AA}}{p_{AA}} - \frac{n_{AT}}{1 - p_{AA} - p_{AC} - p_{AG}}$$

Equate the derivatives to zero. This yields four systems of equations to solve, e.g.:

$$\begin{cases} n_{AT} p_{AA} &= n_{AA} \left(1 - p_{AA} - p_{AC} - p_{AG} \right) \\ n_{AT} p_{AC} &= n_{AC} \left(1 - p_{AA} - p_{AC} - p_{AG} \right) \\ n_{AT} p_{AG} &= n_{AG} \left(1 - p_{AA} - p_{AC} - p_{AG} \right) \end{cases}$$

The transition probabilities are then estimated by:

 n_{AA} $\frac{n_{AA}}{n_{AA} + n_{AC} + n_{AG} + n_{AT}}$ \hat{p}_{AA} n_{AC} $\frac{n_{AC}}{n_{AA} + n_{AC} + n_{AG} + n_{AT}}$ \hat{p}_{AC} n_{TT} \hat{p}_{TT} $n_{TA} + n_{TC} + n_{TG} + n_{TT}$

Verify 2nd order part. derivatives of log-likehood are negative!

Thus, for the following sequence: **ATCGATCGCA**, tabulation yields



The maximum likelihood estimates thus become:

$$\hat{p}_{AA} = 0, \, \hat{p}_{AC} = 0, \, \hat{p}_{AG} = 0, \, \hat{p}_{AT} = 1, \\ \hat{p}_{CA} = \frac{1}{3}, \, \hat{p}_{CC} = 0, \, \hat{p}_{CG} = \frac{2}{3}, \, \hat{p}_{CT} = 0 \\ \dots, \dots, \dots, \dots, \dots$$



- > table(DNAseq)

T 0 2 0 0

DNAseq

- ACGT
- 3 3 2 2
- > table(DNAseq[1:9], DNAseq[2:10])
 A C G T
 A 0 0 0 2
 C 1 0 2 0
 G 1 1 0 0

If one may assume the observed data is a realization of a stationary Markov chain, the initial distribution is estimated by the *stationary distribution* (discussed next time).

If only one realization of a stationary Markov chain is available and stationarity cannot be assumed, the initial distribution is estimated by:

$$\hat{\pi}_i = I_{\{X_1 = E_i\}}$$

Testing the order of the Markov chain

Testing the order of a Markov chain

Often the order of the Markov chain is unknown and needs to be determined from the data. This is done using a *Chi-square test* for independence.

Idea : assess validity of a 0th order Markov chain using data. \rightarrow Under a 0th order Markov chain, e.g.: P(G, T) = P(G) P(T)

- \rightarrow All quantities, P(G, T), P(G) and P(T), can be estimated from the data by their frequencies.
- → Comparison of these frequencies measures the fit of a 0th order Markov chain: obs. # dinucl. (G, T) - exp # dinucl. (G, T).
- \rightarrow Large difference indicate poor fit of 0th order Markov chain.

Testing the order of a Markov chain

Consider a DNA-sequence. To assess the order, one first tests whether the sequence consists of independent letters. Hereto count the nucleotides, di-nucleotides and tri-nucleotides in the sequence:

$$N(x_t) = \{ \#t \mid X_t = x_t \}$$

$$N(x_t, x_{t+1}) = \{ \#t \mid X_t = x_t, X_{t+1} = x_{t+1} \}$$

$$N(x_t, x_{t+1}, x_{t+2}) = \{ \#t \mid X_t = x_t, X_{t+1} = x_{t+1}, X_{t+2} = x_{t+2} \}$$

E.g. for {x₁, x₂, x₃, x₄} = {T, G, A, G}:

$$\rightarrow N(x_t = G) = 2, N(x_t = A) = 1,$$

 $\rightarrow N(x_t = G, x_{t+1} = A) = 1, N(x_t = T, x_{t+1} = A) = 0,$
 $\rightarrow N(x_t = A, x_{t+1} = A, x_{t+2} = A) = 0,$
Assumping independence, we can calculate the expected frequency of a di-nucleotide for a sequence of length *n*:

$$E(\#AG) = (n-1)p(A)p(G)$$
$$= (n-1)\frac{N(A)}{n}\frac{N(G)}{n}$$

To test the order of a Markov chain, we now compare the observed and expected of (say) di-nucleotides.

Compare data from a 1st order Markov chain with what is expected on the basis of a 0th order Markov chain.

Observed dimer frequency

Expected dimer frequency

	Α	C	G	т		Α	C	G	Т
A	98	103	109	65	Α	93.8	133.6	71.0	76.8
C	95	216	37	185	C	133.6	190.2	101.2	109.4
G	81	128	62	13	G	71.0	101.2	53.8	58.2
т	100	87	76	44	т	76.8	109.4	58.2	62.9

According the 0th order Markov model the CG-dimer ought to be observed $\approx 3x$ more than is done.

The null hypothesis of independence between the letters of the DNA sequence evaluated by the χ^2 statistic:

$$\chi^2 = \sum_{x_t, x_{t+1} \in \{A, C, G, T\}} \frac{\left(N(x_t, x_{t+1}) - N(x_t) N(x_{t+1})/(n-1)\right)^2}{N(x_t) N(x_{t+1})/(n-1)}$$

which has $(4-1) \times (4-1) = 9$ degrees of freedom.

If the null hypothesis cannot be rejected, one would fit a Markov model of order m=0. However, if H₀ can be rejected one would test for higher order dependence.



0

To test for 1st order dependence, consider to implications of the 1st order Markov chain model. E.g.:

$$P(X_{t}=G, X_{t+1}=T, X_{t+2}=A) =$$

$$P(X_{t}=G) \times P(X_{t+1}=T | X_{t}=G) \times P(X_{t+2}=A | X_{t+1}=T)$$

Using the definition of conditional probability:

$$P(X_{t}=G, X_{t+1}=T, X_{t+2}=A) =$$

$$P(X_{t+1}=T, X_{t}=G) \times P(X_{t+2}=A, X_{t+1}=T) / P(X_{t}=T).$$

All these probabilities can be estimated directly from the data.

The 1st order dependence hypothesis is evaluated by:

$$\sum_{x_t, x_{t+1}, x_{t+2} \in \{A, C, G, T\}} \frac{\left[N(x_t, x_{t+1}, x_{t+2}) - E\{N(x_t, x_{t+1}, x_{t+2})\}\right]^2}{E(\{N(x_t, x_{t+1}, x_{t+2})\}}$$

where

$$E\{N(x_t, x_{t+1}, x_{t+2})\} = \frac{N(x_t, x_{t+1}) N(x_{t+1}, x_{t+2})}{N(x_{t+1})}$$

This is χ^2 distributed with (16-1) x (4-1) = 45 d.o.f..

A 1st order Markov chain provides a reasonable description of the sequence of the hlyE gene of the E.coli bacteria.

```
Independence test:
  Chi-sq stat: 22.45717, p-value: 0.00754
1<sup>st</sup> order dependence test:
  Chi-sq stat: 55.27470, p-value: 0.14025
```

The sequence of the prrA gene of the E.coli bacteria requires a higher order Markov chain model.

```
Independence test:
  Chi-sq stat: 33.51356, p-value: 1.266532e-04
1<sup>st</sup> order dependence test:
  Chi-sq stat: 114.56290, p-value: 5.506452e-08
```



The independence case, assuming the DNAseq-Object is a character-Object containing the sequence:

- > # calculate nucleotide and dimer frequencies
- > nuclFreq <- matrix(table(DNAseq), ncol=1)</pre>

> dimerFreq <- table(DNAseq[1:(length(DNAseq)-1)], DNAseq[2:length(DNAseq)])

- > # calculate expected dimer frequencies
- > dimerExp <- nuclFreq %*% t(nuclFreq)/(length(DNAseq)-1)</pre>

Exercise: modify the code above for the 1st order test.





We have fitted 1st order Markov chain models to a representative coding and noncoding sequence of the E.coli bacteria.

	P _{coding}					P noncoding				
	Α	C	G	Т		A	C	G	Т	
Α	0.321	0.257	0.211	0.211	A	0.320	0.278	0.231	0.172	
С	0.319	0.207	0.266	0.208	C	0.295	0.205	0.286	0.214	
G	0.259	0.284	0.237	0.219	G	0.241	0.261	0.233	0.265	
т	0.223	0.243	0.309	0.225	Т	0.283	0.238	0.256	0.223	

These models can be used to discriminate between coding and noncoding sequences.

Example: sequence discrimination

For a new sequence with unknown function calculate the likelihood under the coding and noncoding model, e.g.:

$$L_{coding}(\mathbf{X}) = P_{coding}(X_1 = x_1, X_2 = x_2, \dots, X_T = x_T)$$

= $\hat{\pi}_{coding}(x_1) (\hat{\mathbf{P}}_{coding})_{AA}^{n_{AA}} \cdots (\hat{\mathbf{P}}_{coding})_{TT}^{n_{TT}}$

These likelihoods are compared by means of their log ratio: $LR(\mathbf{X}) = \log[L_{coding}(\mathbf{X})] - \log[L_{noncoding}(\mathbf{X})]$ the log-likelihood ratio. If the log likelihood $LR(\mathbf{X})$ of a new sequence exceeds a certain threshold, the sequence is classified as coding and non-coding otherwise.

Back to E.coli

To illustrate the potential of the discrimination approach, we calculate the log likelihood ratio for large set of known coding and noncoding sequences. The distribution of the two sets of LR(X)'s are compared.

Example: sequence discrimination



Conclusion E.coli example

Comparison of the distributions indicate that one could discriminate reasonably between coding and noncoding sequences on the basis of simple 1st order Markov.

Improvements:

- Higher order Markov model.
- Incorporate more structure information of the DNA.



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 : > 2 copies : gain

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.

So far, we only considered one locus. Hence:



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:

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

Using:

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

These probabilities simplify to, e.g.:

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

Parameters α and β determine which aberration type is most prevalent.



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.:

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

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.



Supplementary material: Proof of Kolmogorov-Chapman equations



Proof of the Kolmogorov-Chapman equations:

$$\begin{aligned} (\mathbf{P}^{n+m})_{ij} &= P(X_{n+m} = j \mid X_0 = i) \\ &= \sum_{k=1}^{S} P(X_{n+m} = j, X_n = k \mid X_0 = i) \\ &= \sum_{k=1}^{S} P(X_{n+m} = j, X_n = k, X_0 = i) / P(X_0 = i) \\ &= \sum_{k=1}^{S} P(X_{n+m} = j \mid X_n = k, X_0 = i) P(X_n = k, X_0 = i) / P(X_0 = i) \\ &= \sum_{k=1}^{S} (\mathbf{P}^m)_{kj} P(X_n = k, X_0 = i) / P(X_0 = i) \\ &= \sum_{k=1}^{S} (\mathbf{P}^n)_{ik} (\mathbf{P}^m)_{kj} \end{aligned}$$



SM: Markov processes

Induction proof of $\mathbf{P}^{(n+1)} = \mathbf{P}^{n+1}$ Assume $\mathbf{P}^{(n)} = \mathbf{P}^{n}$

Then:

$$\begin{aligned} (\mathbf{P}^{(n+1)})_{ij} &= P(X_{t+n+1} = E_j \mid X_t = E_i) \\ &= \sum_{k=1}^{S} P(X_{t+n+1} = E_j, X_{t+n} = E_k \mid X_t = E_i) \\ &= \sum_{k=1}^{S} P(X_{t+n+1} = E_j \mid X_{t+n} = E_k) P(X_{t+n} = E_k \mid X_t = E_i) \\ &= \sum_{k=1}^{S} (\mathbf{P})_{kj} (\mathbf{P}^{(n)})_{ik} = \sum_{k=1}^{S} (\mathbf{P})_{kj} (\mathbf{P}^{n})_{ik} = (\mathbf{P}^{n+1})_{ij} \end{aligned}$$

References & further reading

References and further reading

Axelson-Fisk, M. (2010), *Comparative Gene Finding: Models, Algorithms* and Implementation, Springer.

- Basharin, G.P., Langville, A.N., Naumov, V.A (2004), "The life and work of A.A. Markov", *Linear Algebra and its Applications*, **386**, 3-26.
- 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.
- Van Wieringen, W.N., Ros, B.P., Wilting, S.M. (2013), "Modeling the DNA copy number aberration patterns in observational high-throughput cancer data", *Statistical Applications in Genetics and Molecular Biology, 12(2),*143-174.



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

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