

Data Modeling Project for MATH 7241

October 2022

Project

Due date: Friday November 18, week before Thanksgiving.

Groups: if you wish you may collaborate in a group of at most two people. In this case each member of the group must make a substantial contribution to the project, and your project report must describe which contributions were made by each member.

Contents

Step 1: find a good data set!

Recommended source: UCI archive. Search for Time Series.

<https://archive.ics.uci.edu/ml/index.php>



About Citation Policy Donate a Data Set Contact

Repository Web Google Scholar

[View ALL Data Sets](#)

Check out the [beta](#) version of the new UCI Machine Learning Repository we are currently testing! Contact us if you have any issues, questions, or concerns. [Click here to try out the new site.](#)

Welcome to the UC Irvine Machine Learning Repository!

We currently maintain 588 data sets as a service to the machine learning community. You may [view all data sets](#) through our searchable interface. For a general overview of the Repository, please visit our [About page](#). For information about citing data sets in publications, please read our [citation policy](#). If you wish to donate a data set, please consult our [donation policy](#). For any other questions, feel free to [contact the Repository librarians](#).



Latest News:

- 09-24-2018: Welcome to the new Repository admins Dheeru Dua and Efi Karra Taniskidou!
04-04-2013: Welcome to the new Repository admins Kevin Bacho and Moshe Lichman!
03-01-2010: Note from donor regarding Netflix data
10-16-2009: Two new data sets have been added.
09-14-2009: Several new data sets have been added.
03-24-2008: New data sets have been added!
06-25-2007: Two new data sets have been added: UJI Pen Characters, MAGIC Gamma Telescope

Featured Data Set: Iris



Task: Classification
Data Type: Multivariate
Attributes: 4
Instances: 150

Famous database; from Fisher, 1936

Newest Data Sets:

- 04-21-2021:  Synchronous Machine Data Set
04-20-2021:  Wikipedia Math Essentials
04-20-2021:  Wikipedia Math Essentials
02-17-2021:  Hungarian Chickenpox Cases
12-09-2020:  Myocardial infarction complications
10-14-2020:  Gait Classification
10-03-2020:  Codon usage
09-15-2020:  In-vehicle coupon recommendation

Most Popular Data Sets (hits since 2007):

- 4287434:  Iris
2291086:  Adult
1769410:  Wine
1678814:  Wine Quality
1662751:  Heart Disease
1593325:  Breast Cancer Wisconsin (Diagnostic)
1585909:  Bank Marketing
1452321:  Car Evaluation

Search for a time series:

UCI Machine Learning Repository
Center for Machine Learning and Intelligent Systems

Browse Through: 75 Data Sets

	Name	Data Types	Default Task	Attribute Types	# Instances	# Attributes	Year
 Activities of Daily Living (ADLs) Recognition Using Binary Sensors	Multivariate, Sequential, Time-Series	Classification, Clustering			2747		2013
 Activity Recognition from Single Chest-Mounted Accelerometer	Univariate, Sequential, Time-Series	Classification, Clustering	Real				2014
 Activity Recognition system based on Multisensor data fusion (AReM)	Multivariate, Sequential, Time-Series	Classification	Real	42240	6	2016	
 Air Quality	Multivariate, Time-Series	Regression	Real	9358	15	2016	
 Air quality	Multivariate, Time-Series	Regression	Real	9358	15	2016	
 Amazon Access Samples	Time-Series, Domain-Theory	Regression, Clustering, Causal-Discovery		30000	25000	2011	
 Appliances energy prediction	Multivariate, Time-Series	Regression	Real	19756	98	2017	
 Australian Sign Language sign	Multivariate, Time-Series	Classification	Categorical, Real	6650	15	1999	
 Australian Sign Language signs (High Quality)	Multivariate, Time-Series	Classification	Real	2565	22	2002	

[About](#) [Contact](#) [Policy](#) [Donate a Data Set](#) [Logout](#)

Research Web [Google](#)

[View All Data Sets](#)



Select one time series:


Machine Learning Repository
Center for Machine Learning and Intelligent Systems

About Citation Policy Donate a Data Set Contact
 Repository This Google
[View All Data Sets](#)

Air Quality Data Set

[Download](#) [Data Folder](#) [Data Set Description](#)

Abstract: Contains the responses of a gas multisensor device deployed on the field in an Italian city. Hourly responses averages are recorded along with gas concentrations references from a certified analyzer.

Data Set Characteristics:	Multivariate, Time-Series	Number of Instances:	9358	Area:	Computer
Attribute Characteristics:	Real	Number of Attributes:	15	Date Donated:	2016-03-23
Associated Tasks:	Regression	Missing Values?	Yes	Number of Web Hits:	131942

Source:

Severo De Vito (saveto.devito@q3.it), ENIA - National Agency for New Technologies, Energy and Sustainable Economic Development

Data Set Information:

The dataset contains 9358 instances of hourly averaged responses from an array of 5 metal oxide chemical sensors embedded in an Air Quality Chemical Multisensor Device. The device was located on the field in a significantly polluted area, at road level, within an Italian city. Data were recorded from March 2004 to February 2005. The dataset also contains the original hourly available recordings of on-field deployed air quality chemical sensor devices responses. Ground Truth hourly averaged concentrations for CO, Non-Metanic Hydrocarbons, Benzene, Total Nitrogen Oxides (NOx) and Nitrogen Dioxide (NO2) and were provided by a co-located reference certified analyzer. Evidences of cross-sensitivities as well as both concept and sensor drifts are present as described in De Vito et al., *Sens. Actu. A*, Vol. 126, 2, 2006 (citation required) eventually affecting sensors concentration estimation capabilities. Missing values are flagged with -200 value. This dataset can be used exclusively for research purposes. Commercial purposes are fully excluded.

Attribute Information:

- 0 Date (DD/MM/YYYY)
- 1 Time (HH:MM:SS)
- 2 True hourly averaged concentration CO in microg/m³ (reference analyzer)
- 3 True hourly averaged concentration Benzene in microg/m³ (reference analyzer)
- 4 True hourly averaged overall Non-Metanic HydroCarbons concentration in microg/m³ (reference analyzer)
- 5 True hourly averaged Benzene concentration in microg/m³ (reference analyzer)
- 6 PT08-32 (tungsten oxide) hourly averaged sensor response (nominally NMHC targeted)
- 7 PT08-33 (tungsten oxide) hourly averaged sensor response (nominally NO targeted)
- 8 PT08-33 (tungsten oxide) hourly averaged sensor response (nominally NOx targeted)
- 9 True hourly averaged NO2 concentration in microg/m³ (reference analyzer)
- 10 PT08-34 (tungsten oxide) hourly averaged sensor response (nominally NO2 targeted)



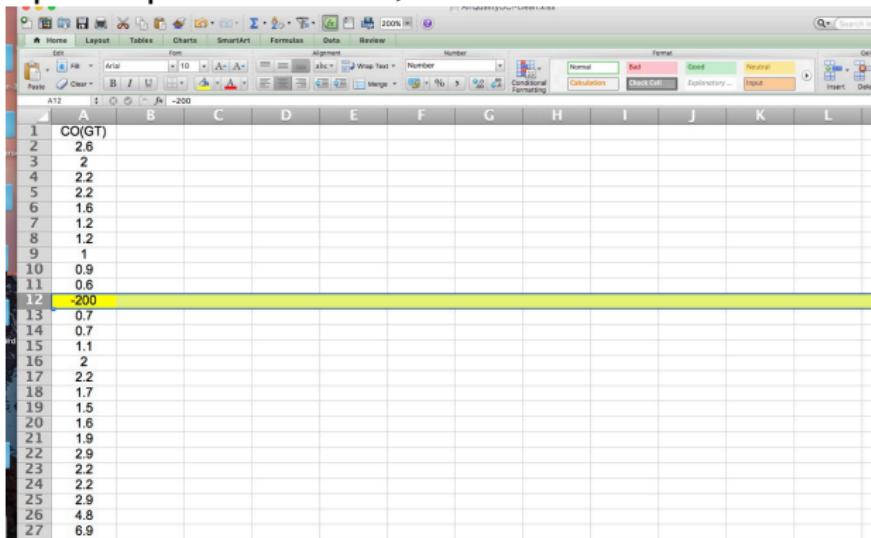
Step 2: download and clean the data set to prepare for modeling. For example, remove unnecessary data, correct 'error' entries etc

Raw data set:

The screenshot shows a Microsoft Excel spreadsheet titled "Project 1 Raw Data Set.xlsx". The data is organized into 10 columns labeled A through J. Column A contains row numbers from 1 to 100. Columns B through J contain various numerical values. The data includes several rows of missing or zero values, such as row 10 where all values are zero, and rows 20, 30, 40, 50, 60, 70, 80, 90, and 100 which also contain mostly zeros. The data appears to represent a time series or a dataset with multiple variables across 100 observations.

	A	B	C	D	E	F	G	H	I	J
1	1	0	0	0	0	0	0	0	0	0
2	31/254	18:00:00	2.5	130	100	100	100	100	100	100
3	31/254	18:00:00	2.5	130	100	100	100	100	100	100
4	31/254	18:00:00	2.5	130	100	100	100	100	100	100
5	31/254	21:00:00	2.5	130	100	100	100	100	100	100
6	31/254	21:00:00	2.5	130	100	100	100	100	100	100
7	31/254	23:00:00	1.2	130	100	100	100	100	100	100
8	31/254	23:00:00	1.2	130	100	100	100	100	100	100
9	31/254	12:00:00	1	130	100	100	100	100	100	100
10	31/254	12:00:00	1	130	100	100	100	100	100	100
11	31/254	12:00:00	1	130	100	100	100	100	100	100
12	31/254	12:00:00	1	130	100	100	100	100	100	100
13	31/254	12:00:00	1	130	100	100	100	100	100	100
14	31/254	12:00:00	1	130	100	100	100	100	100	100
15	31/254	12:00:00	1	130	100	100	100	100	100	100
16	31/254	12:00:00	1	130	100	100	100	100	100	100
17	31/254	12:00:00	1	130	100	100	100	100	100	100
18	31/254	12:00:00	1.5	130	100	100	100	100	100	100
19	31/254	12:00:00	1.5	130	100	100	100	100	100	100
20	31/254	12:00:00	1.5	130	100	100	100	100	100	100
21	31/254	12:00:00	1.5	130	100	100	100	100	100	100
22	31/254	14:00:00	2.5	131	100	100	100	100	100	100
23	31/254	14:00:00	2.5	131	100	100	100	100	100	100
24	31/254	14:00:00	2.5	131	100	100	100	100	100	100
25	31/254	14:00:00	2.5	131	100	100	100	100	100	100
26	31/254	14:00:00	2.5	131	100	100	100	100	100	100
27	31/254	14:00:00	2.5	131	100	100	100	100	100	100
28	31/254	14:00:00	2.5	131	100	100	100	100	100	100
29	31/254	14:00:00	2.5	131	100	100	100	100	100	100
30	31/254	14:00:00	2.5	131	100	100	100	100	100	100
31	31/254	14:00:00	2.5	131	100	100	100	100	100	100
32	31/254	14:00:00	2.5	131	100	100	100	100	100	100
33	31/254	14:00:00	2.5	131	100	100	100	100	100	100
34	31/254	14:00:00	2.5	131	100	100	100	100	100	100
35	31/254	14:00:00	2.5	131	100	100	100	100	100	100
36	31/254	14:00:00	2.5	131	100	100	100	100	100	100
37	31/254	14:00:00	2.5	131	100	100	100	100	100	100
38	31/254	14:00:00	2.5	131	100	100	100	100	100	100
39	31/254	14:00:00	2.5	131	100	100	100	100	100	100
40	31/254	14:00:00	2.5	131	100	100	100	100	100	100
41	31/254	14:00:00	2.5	131	100	100	100	100	100	100
42	31/254	14:00:00	2.5	131	100	100	100	100	100	100
43	31/254	14:00:00	2.5	131	100	100	100	100	100	100
44	31/254	14:00:00	2.5	131	100	100	100	100	100	100
45	31/254	14:00:00	2.5	131	100	100	100	100	100	100
46	31/254	14:00:00	2.5	131	100	100	100	100	100	100
47	31/254	14:00:00	2.5	131	100	100	100	100	100	100
48	31/254	14:00:00	2.5	131	100	100	100	100	100	100
49	31/254	14:00:00	2.5	131	100	100	100	100	100	100
50	31/254	14:00:00	2.5	131	100	100	100	100	100	100
51	31/254	14:00:00	2.5	131	100	100	100	100	100	100
52	31/254	14:00:00	2.5	131	100	100	100	100	100	100
53	31/254	14:00:00	2.5	131	100	100	100	100	100	100
54	31/254	14:00:00	2.5	131	100	100	100	100	100	100
55	31/254	14:00:00	2.5	131	100	100	100	100	100	100
56	31/254	14:00:00	2.5	131	100	100	100	100	100	100
57	31/254	14:00:00	2.5	131	100	100	100	100	100	100
58	31/254	14:00:00	2.5	131	100	100	100	100	100	100
59	31/254	14:00:00	2.5	131	100	100	100	100	100	100
60	31/254	14:00:00	2.5	131	100	100	100	100	100	100
61	31/254	14:00:00	2.5	131	100	100	100	100	100	100
62	31/254	14:00:00	2.5	131	100	100	100	100	100	100
63	31/254	14:00:00	2.5	131	100	100	100	100	100	100
64	31/254	14:00:00	2.5	131	100	100	100	100	100	100
65	31/254	14:00:00	2.5	131	100	100	100	100	100	100
66	31/254	14:00:00	2.5	131	100	100	100	100	100	100
67	31/254	14:00:00	2.5	131	100	100	100	100	100	100
68	31/254	14:00:00	2.5	131	100	100	100	100	100	100
69	31/254	14:00:00	2.5	131	100	100	100	100	100	100
70	31/254	14:00:00	2.5	131	100	100	100	100	100	100
71	31/254	14:00:00	2.5	131	100	100	100	100	100	100
72	31/254	14:00:00	2.5	131	100	100	100	100	100	100
73	31/254	14:00:00	2.5	131	100	100	100	100	100	100
74	31/254	14:00:00	2.5	131	100	100	100	100	100	100
75	31/254	14:00:00	2.5	131	100	100	100	100	100	100
76	31/254	14:00:00	2.5	131	100	100	100	100	100	100
77	31/254	14:00:00	2.5	131	100	100	100	100	100	100
78	31/254	14:00:00	2.5	131	100	100	100	100	100	100
79	31/254	14:00:00	2.5	131	100	100	100	100	100	100
80	31/254	14:00:00	2.5	131	100	100	100	100	100	100
81	31/254	14:00:00	2.5	131	100	100	100	100	100	100
82	31/254	14:00:00	2.5	131	100	100	100	100	100	100
83	31/254	14:00:00	2.5	131	100	100	100	100	100	100
84	31/254	14:00:00	2.5	131	100	100	100	100	100	100
85	31/254	14:00:00	2.5	131	100	100	100	100	100	100
86	31/254	14:00:00	2.5	131	100	100	100	100	100	100
87	31/254	14:00:00	2.5	131	100	100	100	100	100	100
88	31/254	14:00:00	2.5	131	100	100	100	100	100	100
89	31/254	14:00:00	2.5	131	100	100	100	100	100	100
90	31/254	14:00:00	2.5	131	100	100	100	100	100	100
91	31/254	14:00:00	2.5	131	100	100	100	100	100	100
92	31/254	14:00:00	2.5	131	100	100	100	100	100	100
93	31/254	14:00:00	2.5	131	100	100	100	100	100	100
94	31/254	14:00:00	2.5	131	100	100	100	100	100	100
95	31/254	14:00:00	2.5	131	100	100	100	100	100	100
96	31/254	14:00:00	2.5	131	100	100	100	100	100	100
97	31/254	14:00:00	2.5	131	100	100	100	100	100	100
98	31/254	14:00:00	2.5	131	100	100	100	100	100	100
99	31/254	14:00:00	2.5	131	100	100	100	100	100	100
100	31/254	14:00:00	2.5	131	100	100	100	100	100	100

Clean it up: keep one column, remove error entries:



A	B	C	D	E	F	G	H	I	J	K	L
1	CO(GT)										
2	2.6										
3	2										
4	2.2										
5	2.2										
6	1.6										
7	1.2										
8	1.2										
9	1										
10	0.9										
11	0.6										
12	-200										
13	0.7										
14	0.7										
15	1.1										
16	2										
17	2.2										
18	1.7										
19	1.5										
20	1.6										
21	1.9										
22	2.9										
23	2.2										
24	2.2										
25	2.9										
26	4.8										
27	6.9										

Step 3: map your data into a Markov chain. To do this, you must choose the states for your model. Each entry in the time series should map into a unique state.

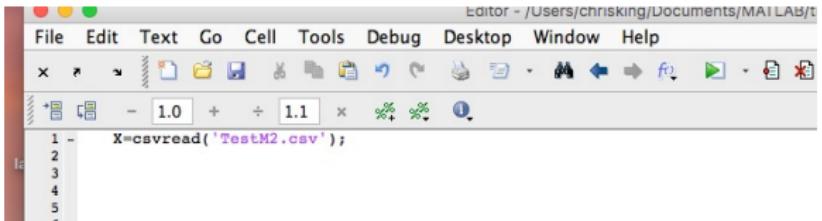
In this example, we choose a 9-state model, with the states $\{1, 2, \dots, 9\}$. Map each entry into a state by rounding to the nearest integer.

AirQualityUCI-clean2.xlsx

	A	B	C	D	E	F	G	H	I
1	CO(GT)		Rounded						
2	2.6		3						
3	2		2						
4	2.2		2						
5	2.2		2						
6	1.6		2						
7	1.2		1						
8	1.2		1						
9	1		1						
10	0.9		1						
11	0.6		1						
12	0.7		1						
13	0.7		1						
14	1.1		1						
15	2		2						
16	2.2		2						
17	1.7		2						
18	1.5		2						
19	1.6		2						
20	1.9		2						
21	2.9		3						
22	2.2		2						
23	2.2		2						
24	2.9		3						
25	4.8		5						
26	6.9		7						

Step 4: transfer the time series to a platform for analysis, eg Excel, Matlab, R etc.

For example, here we import the spreadsheet into Matlab:

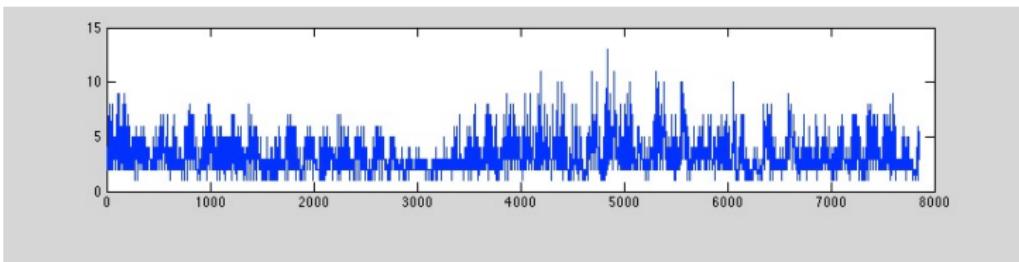


A screenshot of a MATLAB Editor window. The title bar reads "Editor - /Users/chrisking/Documents/MATLAB/t". The menu bar includes File, Edit, Text, Go, Cell, Tools, Debug, Desktop, Window, and Help. Below the menu is a toolbar with various icons. The code area contains the following MATLAB command:

```
X=csvread('TestM2.csv');
```

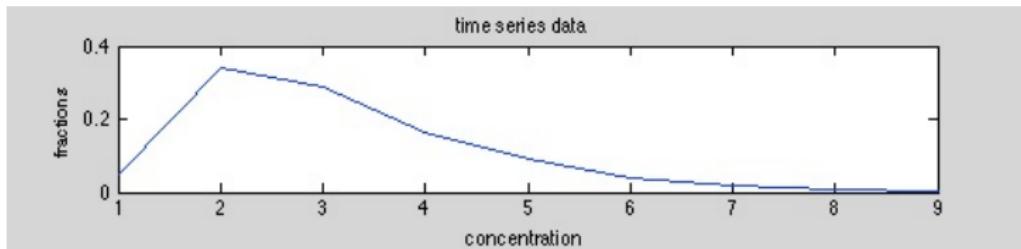
The line numbers 1 through 5 are visible on the left side of the code area.

and here is the complete time series:



Step 5: compute the occupation frequencies for each state, and turn this into a probability distribution. This is the empirical distribution of your chain.

Empirical distribution from time series: fraction of time spent in each state



Step 6: compute the frequencies of jumps between each pair of states. Divide by the occupation frequency at each state so that the total jump probability out of each state is 1. This is your transition matrix.

Transition matrix from time series:

```
Command Window
File Edit Debug Desktop Window Help
New to MATLAB? Watch this Video, see Demos, or read Getting Started.
Trans2 =
Columns 1 through 7
    0.6393    0.3532    0.0050    0.0025         0         0         0
    0.0540    0.7504    0.1507    0.0317    0.0110    0.0015    0.0008
    0.0009    0.2129    0.5973    0.1443    0.0356    0.0058    0.0018
    0         0.0297    0.3143    0.4425    0.1572    0.0414    0.0133
    0         0.0029    0.1225    0.3098    0.3429    0.1614    0.0490
    0         0.0031    0.0404    0.2112    0.2888    0.2360    0.1553
    0         0         0.0132    0.1060    0.2384    0.2914    0.1589
    0         0         0.0161    0.0484    0.2258    0.2097    0.1935
    0         0         0         0.0323    0.0323    0.1935    0.1290
Columns 8 through 9
    0         0
    0         0
    0.0009    0
    0.0016    0
    0.0072    0.0014
    0.0404    0.0248
    0.1258    0.0596
    0.1774    0.0806
    0.1935    0.1935
fx >> |
```

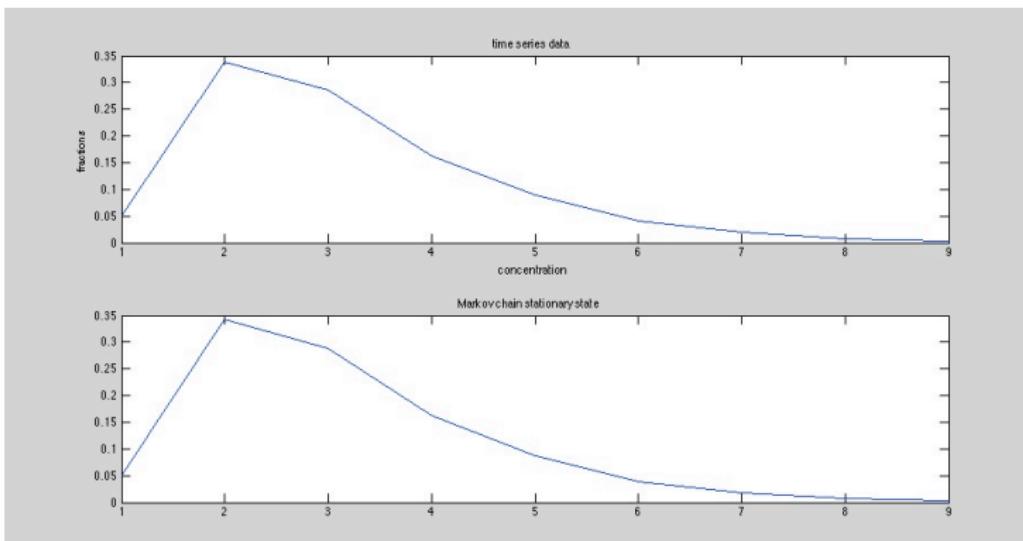
Step 7: find the stationary vector of your transition matrix. In case it is not unique, find all stationary vectors.

Stationary distribution of Markov chain:

```
ans =
0.0522    0.3425    0.2879    0.1623    0.0869    0.0397    0.0180    0.0069    0.0034
fx >>
```

Step 8: compare the empirical distribution of the data set and the stationary distribution of your chain. Note any similarities!

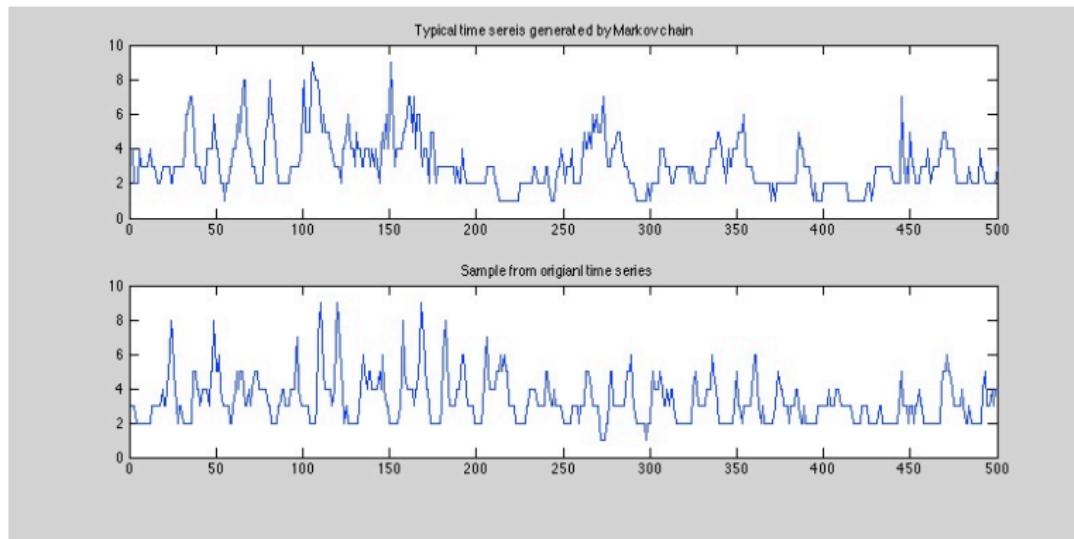
Empirical distribution from time series compared to the stationary distribution of chain:



Step 10:

Build a simulation of the Markov chain, using the transition matrix that you computed in Step 6. Generate a typical time series using your simulation, and compare with the original time series. Does it look like a good model?

Compare simulation of the Markov chain with original time series:



Step 11: Compare your simulation with the original time series using a goodness of fit test for the 2-step transition probabilities. First, working with the original time series, count the number of times the chain jumps from state i to state j in *two successive steps*. Let N_{ij} denote this number, for each pair i, j (note that you should include the cases where $i = j$). Also define

$$N_i = \sum_j N_{ij} \quad \text{for each state } i.$$

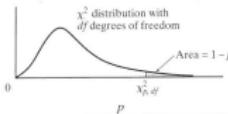
Second, let \hat{p}_{ij} be the transition matrix that you computed in Step 6 above. Define

$$M_{ij} = N_i \sum_k \hat{p}_{ik} \hat{p}_{kj}$$

Step 11 (cont.): Use a goodness of fit test to compare the *observed frequencies* $\{N_{ij}\}$ with the *expected frequencies* $\{M_{ij}\}$ at the 5% significance level, for all pairs i, j . Note: if $M_{ij} < 4$ for some pair i, j , then you should pool frequencies for different values of j (for the same i) in order to get the frequency above 4 for each entry.

Based on your results, decide if the 2-step Markov chain is a good model for the 2-step transitions of the original data set.

Step 12: write a report (maximum 7 pages) explaining how you carried out the above steps, including: source and nature of raw data set, how it was cleaned, choice of states for Markov chain, empirical distribution, transition matrix, stationary distribution, comparison of empirical and stationary distributions, goodness of fit test. At the end of your report, answer this question: 'do you consider that the Markov chain method produces a good model for this time series? Explain your answer'.

Table A.3 Upper and Lower Percentiles of χ^2 Distributions

df	0.010	0.025	0.050	0.10	0.90	0.95	0.975	0.99
1	0.000157	0.000982	0.00393	0.0158	2.708	3.841	5.024	6.615
2	0.0201	0.0506	0.103	0.211	4.605	5.991	7.378	9.210
3	0.115	0.216	0.352	0.584	6.251	7.815	9.348	11.345
4	0.297	0.484	0.711	1.064	7.779	9.488	11.143	13.277
5	0.554	0.831	1.145	1.610	9.236	11.070	12.832	15.086
6	0.872	1.237	1.635	2.304	10.645	12.592	14.449	16.812
7	1.239	1.690	2.167	2.833	12.017	14.067	16.013	18.475
8	1.646	2.180	2.733	3.490	13.362	15.307	17.535	20.090
9	2.088	2.700	3.325	4.168	14.684	16.919	19.023	21.666
10	2.558	3.247	3.940	4.865	15.987	18.307	20.483	23.209
11	3.053	3.816	4.575	5.578	17.275	19.675	21.920	24.725
12	3.571	4.404	5.226	6.304	18.549	21.026	23.330	26.217
13	4.107	5.009	5.892	7.042	19.812	22.362	24.736	27.688
14	4.660	5.629	6.571	7.790	21.064	23.685	26.119	29.141
15	5.229	6.262	7.261	8.547	22.307	24.996	27.488	30.578
16	5.812	6.908	7.962	9.312	23.542	26.296	28.845	32.000
17	6.408	7.564	8.672	10.085	24.769	27.587	30.191	33.409
18	7.015	8.231	9.396	10.865	25.989	28.869	31.525	34.805
19	7.633	8.907	10.117	11.651	27.204	30.144	32.852	36.191
20	8.260	9.591	10.851	12.443	28.412	31.410	34.170	37.566
21	8.897	10.283	11.591	13.240	29.615	32.671	35.479	38.932
22	9.542	10.982	12.338	14.041	30.813	33.924	36.781	40.289
23	10.196	11.688	13.091	14.848	32.007	35.172	38.076	41.638
24	10.856	12.401	13.848	15.659	33.196	36.415	39.364	42.980
25	11.524	13.120	14.611	16.473	34.382	37.652	40.646	44.314
26	12.194	13.844	15.379	17.292	35.563	38.885	41.923	45.642
27	12.879	14.573	16.151	18.114	36.741	40.113	43.194	46.963
28	13.565	15.308	16.926	18.939	37.916	41.337	44.461	48.278
29	14.256	16.047	17.708	19.768	39.087	42.557	45.722	49.588
30	14.953	16.791	18.493	20.599	40.256	43.773	46.979	50.892
31	15.655	17.539	19.281	21.434	41.422	44.985	48.232	52.191
32	16.362	18.291	20.072	22.271	42.585	46.194	49.480	53.486
33	17.073	19.047	20.867	23.110	43.745	47.400	50.725	54.776
34	17.789	19.806	21.664	23.952	44.903	48.602	51.966	56.061

0.4 Absorbing chains

Definition 4 A state i is absorbing if $p_{ii} = 1$. A chain is absorbing if for every state i there is an absorbing state which is accessible from i . A non-absorbing state in an absorbing chain is called a transient state.

Consider an absorbing chain with r absorbing states and t transient states. Re-order the states so that the transient states come first, then the absorbing states. The transition matrix then has the form

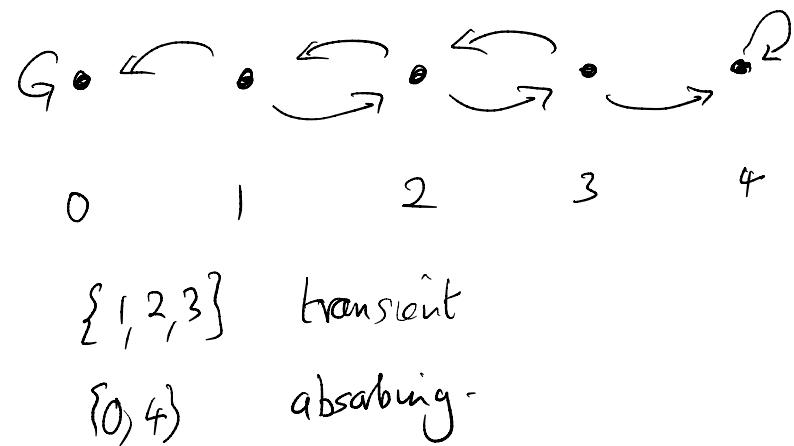
$$P = \begin{pmatrix} Q & R \\ 0 & I \end{pmatrix} \quad (77)$$

where I is the $r \times r$ identity matrix.

$$P = \begin{pmatrix} Q & | & R \\ \hline 0 & | & I \end{pmatrix} \begin{array}{l} \text{transient states} \\ \text{absorbing states} \end{array}$$

Example 7 For the drunkard's walk, show that

$$Q = \begin{pmatrix} 0 & 1/2 & 0 \\ 1/2 & 0 & 1/2 \\ 0 & 1/2 & 0 \end{pmatrix}, \quad R = \begin{pmatrix} 1/2 & 0 \\ 0 & 0 \\ 0 & 1/2 \end{pmatrix}, \quad I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad (78)$$



Simple calculations show that for all $n \geq 1$

$$P^n = \begin{pmatrix} Q^n & R_n \\ 0 & T \end{pmatrix} \quad (79)$$

where R_n is a complicated matrix depending on Q and R .

Lemma 2 As $n \rightarrow \infty$,

~~transient~~
for all ~~absorbing~~ states i, j . $(Q^n)_{ij} \rightarrow 0$ [See Ex. 5 on Problem Set 5]

Proof: for a transient state i , there is an absorbing state k , an integer n_i and $\delta_i > 0$ such that

$$p_{ik}(n_i) = \delta_i > 0 \quad (80)$$

Let $n = \max n_i$, and $\delta = \min \delta_i$, then for any $i \in T$, there is a state $k \in R$ such that

$$p_{ik}(n) \geq \delta \quad (81)$$

Hence for any $i \in T$,

$$\sum_{j \in T} Q_{ij}^n = 1 - \sum_{k \in R} P_{ik}^n = 1 - \sum_{k \in R} p_{ik}(n) \leq 1 - \delta \quad (82)$$

In particular this means that $Q_{ij}^n \leq 1 - \delta$ for all $i, j \in T$. So for all $i \in T$ we get

$$\sum_{j \in T} Q_{ij}^{2n} = \sum_{k \in T} Q_{ik}^n \sum_{j \in T} Q_{kj}^n \leq (1 - \delta) \sum_{k \in T} Q_{ik}^n \leq (1 - \delta)^2 \quad (83)$$

This iterates to give

$$\sum_{j \in T} Q_{ij}^{kn} \rightarrow 0 \quad \text{as } k \rightarrow \infty \quad (84)$$

for all $i \in T$.

It remains to notice that

$$\sum_{j \in T} Q_{ij}^{m+1} = \sum_{k \in T} Q_{ik}^m \sum_{j \in T} Q_{kj}^m \leq \sum_{k \in T} Q_{ik}^m \quad (85)$$

and hence the sequence $\{\sum_{k \in T} Q_{ik}^m\}$ is monotone decreasing in m . Therefore

$$\sum_{j \in T} Q_{ij}^k \rightarrow 0 \quad \text{as } k \rightarrow \infty \quad (86)$$

for all $i \in T$, which proves the result.

QED

Notice what the result says: the probability of remaining in the transient states goes to zero, so eventually the chain must transition to the absorbing states. So the quantities of interest are related to the time (=number of steps) needed until the chain exits the transient states and enters the absorbing states, and the number of visits to other transient states.

Consider the equation

$$x = Qx \quad (87)$$

Applying Q to both sides we deduce that

$$x = Q^2x \quad (88)$$

and iterating this leads to

$$x = Q^n x \quad (89)$$

for all n . Since $Q^n \rightarrow 0$ it follows that $x = 0$. Hence there is no nonzero solution of the equation $x = Qx$ and therefore the matrix $I - Q$ is non-singular and so invertible.

Define the fundamental matrix

$$N = (I - Q)^{-1} \quad (90)$$

Note that

$$(I + Q + Q^2 + \cdots + Q^n)(I - Q) = I - Q^{n+1} \quad (91)$$

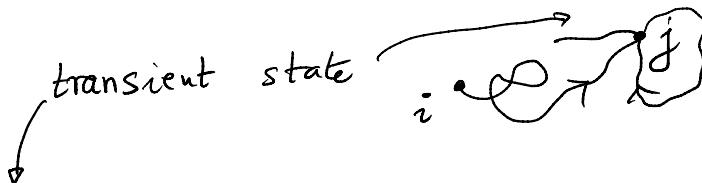
and letting $n \rightarrow \infty$ we deduce that

$$N = I + Q + Q^2 + \cdots \quad (92)$$

Theorem 5 Let i, j be transient states. Then

- (1) N_{ij} is the expected number of visits to state j starting from state i (counting initial state if $i = j$).
- (2) $\sum_j N_{ij}$ is the expected number of steps of the chain, starting in state i , until it is absorbed.
- (3) define the $t \times r$ matrix $B = NR$. Then B_{ik} is the probability that the chain is absorbed in state k , given that it started in state i .

$$\begin{aligned} P &= \begin{pmatrix} Q & R \\ 0 & I \end{pmatrix} & N &= (I - Q)^{-1} \\ B &= NR \end{aligned}$$



Proof: the chain starts at $X_0 = i$. Given a state $j \in T$, for $k \geq 0$ define indicator random variables as follows:

$$Y^{(k)} = \begin{cases} 1 & \text{if } X_k = j \\ 0 & \text{else} \end{cases} \quad \text{indicator r.v.} \quad (93)$$

Then for $k \geq 1$

$$\mathbb{E}Y^{(k)} = P(Y^{(k)} = 1) = P(X_k = j) = p_{ij}(k) = (Q^k)_{ij} \quad (94)$$

and for $k = 0$ we get $\mathbb{E}Y^{(0)} = \delta_{ij}$.

Now the number of visits to the state j in the first n steps is $Y^{(0)} + Y^{(1)} + \dots + Y^{(n)}$. Taking the expected value yields the sum

$$E[\# \text{ visits to } j] = \delta_{ij} + Q_{ij} + (Q^2)_{ij} + \dots + (Q^n)_{ij} = (I + Q + Q^2 + \dots + Q^n)_{ij} \rightarrow N_{ij} \quad (95)$$

as $n \rightarrow \infty$.

which converges to N_{ij} as $n \rightarrow \infty$. This proves (1).

For (2), note that the sum of visits to all transient states is the total number of steps of the chain before it leaves the transient states. For (3), use $N = \sum Q^n$ to write

$$\begin{aligned} B_{ik} &= (NR)_{ik} = \sum_{j \in T} N_{ij} R_{jk} \\ &= \sum_{j \in T} \sum_{n=0}^{\infty} (Q^n)_{ij} R_{jk} \\ &= \sum_{n=0}^{\infty} \sum_{j \in T} (Q^n)_{ij} R_{jk} \end{aligned} \quad \text{prob. } i \xrightarrow{n+1} k \text{ in } (n+1) \text{ steps.} \quad (96)$$

and note that $\sum_{j \in T} (Q^n)_{ij} R_{jk}$ is the probability that the chain takes n steps to transient states before exiting to the absorbing state k . Since this is the only way that the chain can transition to k in $n+1$ steps, the result follows.

QED

Example 8 For the drunkard's walk,

$$Q^{2n+1} = 2^{-n}Q, \quad Q^{2n+2} = 2^{-n}Q^2 \quad (97)$$

and

$$N = \begin{pmatrix} 3/2 & 1 & 1/2 \\ 1 & 2 & 1 \\ 1/2 & 1 & 3/2 \end{pmatrix} \quad (98)$$

Also

$$B = NR = \begin{pmatrix} 3/4 & 1/4 \\ 1/2 & 1/2 \\ 1/4 & 3/4 \end{pmatrix} \quad (99)$$

$$Q = \begin{pmatrix} 0 & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & \frac{1}{2} & 0 \end{pmatrix} \quad R = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & 0 \\ 0 & \frac{1}{2} \end{pmatrix}$$

$$I - Q = \begin{pmatrix} 1 & -\frac{1}{2} & 0 \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ 0 & -\frac{1}{2} & 1 \end{pmatrix} \quad \text{all positive entries!}$$

$$N = (I - Q)^{-1} = \begin{pmatrix} \frac{3}{2} & 1 & \frac{1}{2} \\ 1 & 2 & 1 \\ \frac{1}{2} & 1 & \frac{3}{2} \end{pmatrix} \quad \begin{matrix} 1 \\ 2 \\ 3 \end{matrix}$$

- start in state 1, how many visits to state 3 before absorbed?

$$\text{mean \# visits } 1 \rightarrow 3 = N_{13} = \frac{1}{2}$$

- start in state 2, how many steps until the chain is absorbed?

$$\text{mean \# steps} = N_{21} + N_{22} + N_{23}$$

$$= 4.$$

$$B = NR = \begin{pmatrix} \frac{3}{2} & 1 & \frac{1}{2} \\ 1 & 2 & 1 \\ \frac{1}{2} & 1 & \frac{3}{2} \end{pmatrix} \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & 0 \\ 0 & \frac{1}{2} \end{pmatrix}$$

$3 \times 3 \qquad \qquad \qquad 3 \times 2$

$$= \begin{pmatrix} \frac{3}{4} & \boxed{\frac{1}{4}} & 1 \\ \frac{1}{2} & \frac{1}{2} & 2 \\ \frac{1}{4} & \frac{3}{4} & 3 \\ 0 & 4 & \end{pmatrix}$$

- start in state 1, what is prob.
to be absorbed in state 4?

$$B_{14} = \frac{1}{4}$$

0.5 Addendum on transient states

There are many interesting problems where transient states are important. We will not say too much about them, but just note that there are generally two questions of interest:

- How long until the chain exits the transient states?
- Which irreducible class will it then enter?

Both questions can be answered in a systematic way by using the *associated absorbing chain* together with its *fundamental matrix*.

Suppose the original chain is decomposed as

$$\Omega = T \cup C_1 \cup C_2 \cup \dots$$

where T denotes the transient states, and each C_i is a closed irreducible class.

Then the associated absorbing chain has state space

$$\Omega^{abs} = T \cup \{c_1\} \cup \{c_2\} \cup \dots$$

where T are the same transient states and each closed irreducible class C_i has been replaced by a *single state* c_i .

With respect to this decomposition the transition matrix of the chain on Ω^{abs} is

$$P^{abs} = \begin{pmatrix} Q & R \\ 0 & \mathbb{I} \end{pmatrix}$$

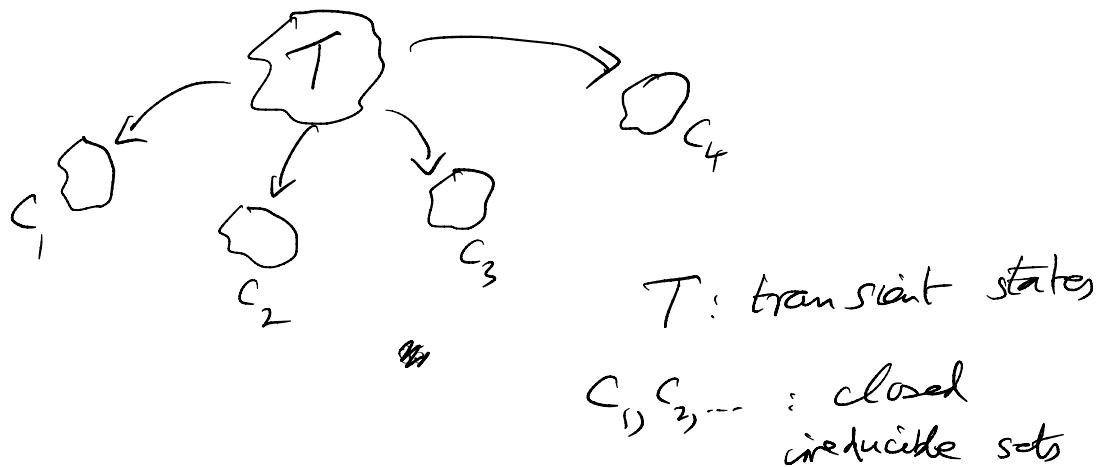
where Q is the original transition matrix between transient states, and

$$R_{tj} = \sum_{a \in C_j} p_{ta}$$

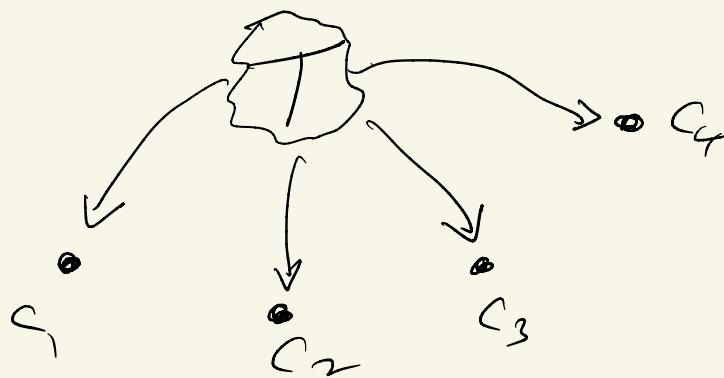
where t is a transient state and the sum runs over all states in the class C_j . Finally \mathbb{I} is the identity matrix on the states $\{c_1\} \cup \{c_2\} \cup \dots$

Then the fundamental matrix of the absorbing chain is defined to be

$$N = (\mathbb{I} - Q)^{-1}$$



Question: starting in T , how long until enter a closed class? what are the prob's for classes?



Build new chain: replace each class C_i by a single absorbing state c_i .

$$P = \begin{pmatrix} Q & R \\ - & I \end{pmatrix}^T \text{Abs.}$$

$Q_{ij} = P_{ij}$ same as
rigid chain.

$$R_{ik} = \sum_{a \in C_k} P_{ia}$$

$$N = (I - Q)^{-1}$$

N_{ij}, B_{ik} same meaning
as before.

Lemma 3 • The matrix element N_{ij} is the expected number of visits of chain to the transient state j starting in the transient state i .

- $\sum_j N_{ij}$ is the expected number of steps until the chain exits the transient states, starting in the transient state i .
- $\sum_k N_{ik} R_{kj}$ is the probability that the chain will enter the class C_j starting in the transient state i .

This works well for a small number of transient states, but becomes cumbersome for larger numbers.

Example Gambler's Ruin 



X_k = gambler's fortune after k steps.

$$\text{fair game} \Rightarrow X_{k+1} = \begin{cases} X_k + 1 & \text{prob} = \frac{1}{2} \\ X_k - 1 & \text{prob} = \frac{1}{2}. \end{cases}$$

Simple random walk on $\{0, 1, 2, \dots, N\}$.

0 = absorbing state = ruin.

N = absorbing state = win.

P_k = prob. that³⁶ gambler wins
starting at $X_0 = k$.

$$P_k = \mathbb{P}(X_n = N \text{ some } n \geq 1 \mid X_0 = k)$$

= prob. to absorb at state N .

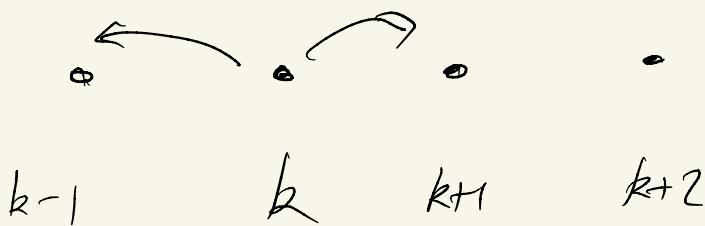
Absorbing chain!

$$P = \begin{pmatrix} Q & R \\ 0 & I \end{pmatrix}$$

Answer: $B_{kN} = (NR)_{kN}$.

State space is too large.

Calculate P_k by conditioning:



$$X_0 = k,$$

Condition on first step:

$$\text{(X)} \quad P_k = \frac{1}{2} P_{k+1} + \frac{1}{2} P_{k-1}$$

$$k=1, 2, \dots, N-1$$

Boundary condition:

$$P_0 = 0, \quad P_N = 1.$$

$$\text{(X)} \quad \frac{1}{2} P_k + \frac{1}{2} P_k = \frac{1}{2} P_{k+1} + \frac{1}{2} P_{k-1}$$

$$\frac{1}{2} P_k - \frac{1}{2} P_{k-1} = \frac{1}{2} P_{k+1} - \frac{1}{2} P_k$$

$$P_k - P_{k-1} = P_{k+1} - P_k.$$

Let $U_k = P_k - P_{k-1}$

$$(k=1, 2, \dots, N).$$

$$\boxed{U_k = U_{k+1}}$$

$$k=1, 2, \dots, N-1$$

$$U_2 = U_1$$

$$U_3 = U_1$$

$$U_4 = U_1$$

$$\boxed{U_k = U_1}$$

$$k=1, 2, \dots, N-1.$$

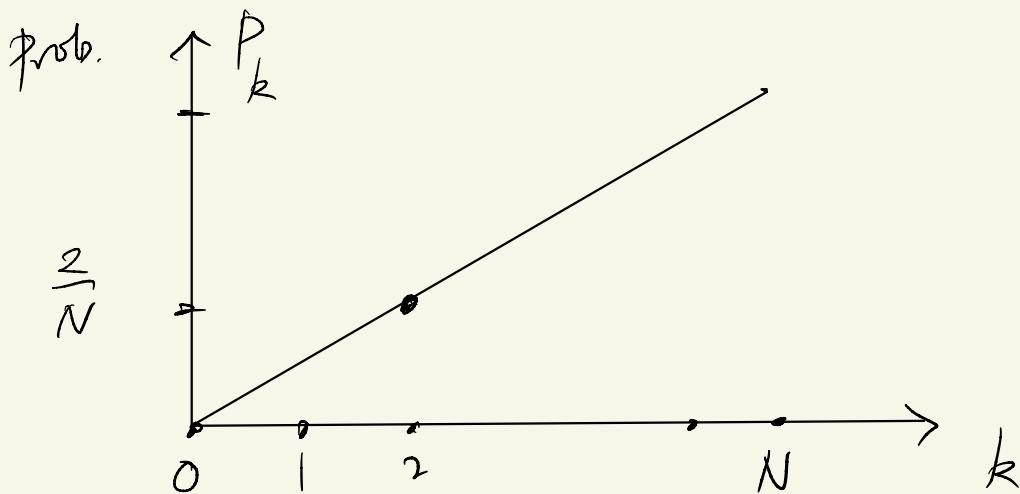
$$u_k = p_k - p_{k-1}$$

$$p_k = k u_1 + b.$$

B.C. $p_0 = 0 \Rightarrow b = 0$

$$p_N = 1 \Rightarrow 1 = N u_1$$

$$\Rightarrow \boxed{p_k = \frac{k}{N}} \quad k=0, 1, 2, \dots, N.$$



0.6 Convergence to the stationary distribution.

An irreducible persistent chain has a unique stationary distribution. It is reasonable to expect that

$$P(X_n = k | X_0 = i) \rightarrow \pi_k$$

Perron–Frobenius for regular chain.

as $n \rightarrow \infty$. We determine the cases where this holds.

The extra condition needed concerns the period of the chain. For example, the RW has period 2: we saw that $P(X_n = 0 | X_0 = 0) = 0$ unless n is even. This is a special case of periodic behavior.

Definition 5 *The period of the state i is*

$$d(i) = \gcd\{n \mid p_{ii}(n) > 0\}$$

The state i is aperiodic if $d(i) = 1$.

Definition 6 *An irreducible, aperiodic, persistent Markov chain is called ergodic.*



Theorem 6 For an ergodic chain,

$$p_{ij}(n) \rightarrow \pi_j = \frac{1}{\mu_j} \quad (100)$$

as $n \rightarrow \infty$, for all $i, j \in S$.

converges as $n \rightarrow \infty$

π_j = stationary prob.
for state j

μ_j = mean first return
time to state j .

0.7 Mixing time for ergodic chain

The ergodic Theorem says that $P(X_n = j) \rightarrow \pi_j$ as $n \rightarrow \infty$ for every state j . So the chain ‘forgets’ about its initial state and settles into the stationary distribution.

How long does this process take? This is an interesting question for many applications.

Card shuffling

Persi Diaconis (b. 1945) left school at 14 to travel with a conjurer and learn card tricks. After returning to college he supported himself by playing poker on ships between New York and South America. Diaconis is currently a Professor at Stanford University.

Diaconis famously showed that it takes seven shuffles to randomize a deck of cards. There is a precise statement of this, but roughly it means that starting from an ordered deck, and repeatedly using the ‘riffle shuffle’, the distance from the uniform distribution (which is the stationary distribution in this case) decreases and gets below $1/2$ after 7 shuffles. After this it drops by a factor $1/2$ at each subsequent shuffle.

In general convergence to the stationary distribution occurs at an exponential rate (determined by the eigenvalues of the transition matrix).

Rate of convergence

[Seneta]: another way to express the Perron-Frobenius result is to say that for the matrix P , 1 is the largest eigenvalue (in absolute value) and w is the unique eigenvector (up to scalar multiples). Let λ_2 be the second largest eigenvalue of P so that $1 > |\lambda_2| \geq |\lambda_i|$. Let m_2 be the multiplicity of λ_2 . Then the following estimate holds: there is $C < \infty$ such that for all $n \geq 1$

$$\|P^n - ew^T\| \leq C n^{m_2-1} |\lambda_2|^n \quad (101)$$

So the convergence $P^n \rightarrow ew^T$ is exponential with rate determined by the first spectral gap.

$$e w^T = \begin{pmatrix} w_1 & w_2 & \cdots & w_m \\ w_1 & w_2 & \cdots & w_m \\ \vdots & \vdots & \ddots & \vdots \end{pmatrix}$$

$w^T P = w^T$: left eigenvector with eigenvalue 1.