

# COMS 4771 Lecture 22

## 1. Markov models

# MARKOV MODELS

# SEQUENCE MODELS

A **sequence model** (or **time series model**) is a family of probability distributions for (possibly infinite) *sequences* of random variables  $\{X_t\}_{t \in \mathcal{T}}$ .

# SEQUENCE MODELS

A **sequence model** (or **time series model**) is a family of probability distributions for (possibly infinite) *sequences* of random variables  $\{X_t\}_{t \in \mathcal{T}}$ .

- ▶  $\{X_t\}_{t \in \mathcal{T}}$  is a **stochastic process** indexed by the totally-ordered set  $\mathcal{T}$  (e.g.,  $\mathcal{T} = \mathbb{N}$  for discrete time series).

# SEQUENCE MODELS

A **sequence model** (or **time series model**) is a family of probability distributions for (possibly infinite) *sequences* of random variables  $\{X_t\}_{t \in \mathcal{T}}$ .

- ▶  $\{X_t\}_{t \in \mathcal{T}}$  is a **stochastic process** indexed by the totally-ordered set  $\mathcal{T}$  (e.g.,  $\mathcal{T} = \mathbb{N}$  for discrete time series).
- ▶ Special emphasis is placed on the **linear ordering of  $\mathcal{T}$** .

# SEQUENCE MODELS

A **sequence model** (or **time series model**) is a family of probability distributions for (possibly infinite) *sequences* of random variables  $\{X_t\}_{t \in \mathcal{T}}$ .

- ▶  $\{X_t\}_{t \in \mathcal{T}}$  is a **stochastic process** indexed by the totally-ordered set  $\mathcal{T}$  (e.g.,  $\mathcal{T} = \mathbb{N}$  for discrete time series).
- ▶ Special emphasis is placed on the **linear ordering of  $\mathcal{T}$** .

If  $t \in \mathcal{T}$  is the “current time”, then  $X_t$  is the “current state”;  
 $X_\tau$  for  $\tau < t$  are “past states”; and  $X_\tau$  for  $\tau > t$  are “future states”.

# SEQUENCE MODELS

A **sequence model** (or **time series model**) is a family of probability distributions for (possibly infinite) *sequences* of random variables  $\{X_t\}_{t \in \mathcal{T}}$ .

- ▶  $\{X_t\}_{t \in \mathcal{T}}$  is a **stochastic process** indexed by the totally-ordered set  $\mathcal{T}$  (e.g.,  $\mathcal{T} = \mathbb{N}$  for discrete time series).
- ▶ Special emphasis is placed on the **linear ordering of  $\mathcal{T}$** .

If  $t \in \mathcal{T}$  is the “current time”, then  $X_t$  is the “current state”;  $X_\tau$  for  $\tau < t$  are “past states”; and  $X_\tau$  for  $\tau > t$  are “future states”.

(May interchange “state” and “observation”—no distinction for now.)

# SEQUENCE MODELS

A **sequence model** (or **time series model**) is a family of probability distributions for (possibly infinite) *sequences* of random variables  $\{X_t\}_{t \in \mathcal{T}}$ .

- ▶  $\{X_t\}_{t \in \mathcal{T}}$  is a **stochastic process** indexed by the totally-ordered set  $\mathcal{T}$  (e.g.,  $\mathcal{T} = \mathbb{N}$  for discrete time series).
- ▶ Special emphasis is placed on the **linear ordering of  $\mathcal{T}$** .

If  $t \in \mathcal{T}$  is the “current time”, then  $X_t$  is the “current state”;  $X_\tau$  for  $\tau < t$  are “past states”; and  $X_\tau$  for  $\tau > t$  are “future states”.

(May interchange “state” and “observation”—no distinction for now.)

Sequence / time series modeling is an entire subfield in statistics, largely due to the plethora of sequence / time series data in applications:

- ▶ Economic / financial data over time
- ▶ Climate science
- ▶ Genomic sequences
- ▶ Speech and natural language
- ▶ ...



# MARKOV MODELS

A stochastic process  $\{X_t\}_{t \in \mathbb{N}}$  has the **Markov property** if the conditional distribution of the next state  $X_{t+1}$  given all previous states  $\{X_\tau : \tau \leq t\}$  only depends on the value of the current state  $X_t$ .

# MARKOV MODELS

A stochastic process  $\{X_t\}_{t \in \mathbb{N}}$  has the **Markov property** if the conditional distribution of the next state  $X_{t+1}$  given all previous states  $\{X_\tau : \tau \leq t\}$  only depends on the value of the current state  $X_t$ .

If the  $X_t$  are discrete-valued, then the Markov property means that

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

$$\cdots \longrightarrow X_{t-1} \longrightarrow X_t \longrightarrow X_{t+1} \longrightarrow \cdots$$

# MARKOV MODELS

A stochastic process  $\{X_t\}_{t \in \mathbb{N}}$  has the **Markov property** if the conditional distribution of the next state  $X_{t+1}$  given all previous states  $\{X_\tau : \tau \leq t\}$  only depends on the value of the current state  $X_t$ .

If the  $X_t$  are discrete-valued, then the Markov property means that

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

$$\cdots \longrightarrow X_{t-1} \longrightarrow X_t \longrightarrow X_{t+1} \longrightarrow \cdots$$

A stochastic process with the Markov property is called a **Markov chain**.

# MARKOV MODELS

A stochastic process  $\{X_t\}_{t \in \mathbb{N}}$  has the **Markov property** if the conditional distribution of the next state  $X_{t+1}$  given all previous states  $\{X_\tau : \tau \leq t\}$  only depends on the value of the current state  $X_t$ .

If the  $X_t$  are discrete-valued, then the Markov property means that

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

$$\cdots \longrightarrow X_{t-1} \longrightarrow X_t \longrightarrow X_{t+1} \longrightarrow \cdots$$

A stochastic process with the Markov property is called a **Markov chain**.

A sequence model for a Markov chain is called a **Markov model**.

# MARKOV CHAIN DISTRIBUTIONS

To specify a Markov chain (MC):

- ▶ Specify the distribution of the initial state  $X_1$ .
- ▶ Specify a **transition kernel**:  $\Pr(X_{t+1} = x' \mid X_t = x)$  for all  $(x, x')$ .

(Nothing to do with *kernels* as in SVMs/kernel trick/RKHS.)

# MARKOV CHAIN DISTRIBUTIONS

To specify a Markov chain (MC):

- ▶ Specify the distribution of the initial state  $X_1$ .
- ▶ Specify a **transition kernel**:  $\Pr(X_{t+1} = x' \mid X_t = x)$  for all  $(x, x')$ .

(Nothing to do with *kernels* as in SVMs/kernel trick/RKHS.)

We focus on MCs where the **state space** (possible values for each  $X_t$ ) is finite.

For simplicity, we'll assume the state space is  $[d] := \{1, 2, \dots, d\}$ .

# MARKOV CHAIN DISTRIBUTIONS

To specify a Markov chain (MC):

- ▶ Specify the distribution of the initial state  $X_1$ .
- ▶ Specify a **transition kernel**:  $\Pr(X_{t+1} = x' \mid X_t = x)$  for all  $(x, x')$ .

(Nothing to do with *kernels* as in SVMs/kernel trick/RKHS.)

We focus on MCs where the **state space** (possible values for each  $X_t$ ) is finite.

For simplicity, we'll assume the state space is  $[d] := \{1, 2, \dots, d\}$ .

- ▶ Initial state distribution given by a  $d$ -dimensional probability vector  $\pi$

$$\pi_i = \Pr(X_1 = i).$$

# MARKOV CHAIN DISTRIBUTIONS

To specify a Markov chain (MC):

- ▶ Specify the distribution of the initial state  $X_1$ .
- ▶ Specify a **transition kernel**:  $\Pr(X_{t+1} = x' \mid X_t = x)$  for all  $(x, x')$ .

(Nothing to do with *kernels* as in SVMs/kernel trick/RKHS.)

We focus on MCs where the **state space** (possible values for each  $X_t$ ) is finite.

For simplicity, we'll assume the state space is  $[d] := \{1, 2, \dots, d\}$ .

- ▶ Initial state distribution given by a  $d$ -dimensional probability vector  $\pi$

$$\pi_i = \Pr(X_1 = i).$$

- ▶ Transition kernel can be written as a  $d \times d$  matrix  $\mathbf{A}$

$$A_{i,j} = \Pr(X_{t+1} = j \mid X_t = i)$$

Great!!!

(rows of  $\mathbf{A}$  are probability vectors).

Also called a **transition matrix** or **(right) stochastic matrix**.

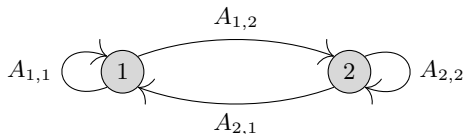


# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\boldsymbol{\pi} = \begin{matrix} & \text{state 1} \\ \text{state 1} & 0.1 \\ \text{state 2} & 0.9 \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & 0.3 & 0.7 \\ \text{state 2} & 0.6 & 0.4 \end{matrix}.$$

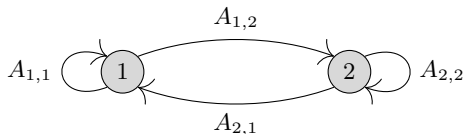


# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\boldsymbol{\pi} = \begin{matrix} & \text{state 1} \\ \text{state 1} & 0.1 \\ \text{state 2} & 0.9 \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & 0.3 & 0.7 \\ \text{state 2} & 0.6 & 0.4 \end{matrix}.$$



A random state sequence drawn from this MC:

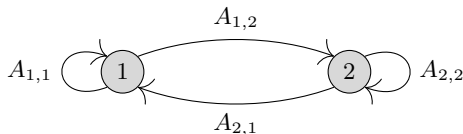
(2, 2, 2, 1, 1, 2, 2, 1)

# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\pi = \begin{matrix} & \text{state 1} \\ \text{state 1} & 0.1 \\ \text{state 2} & 0.9 \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & 0.3 & 0.7 \\ \text{state 2} & 0.6 & 0.4 \end{matrix}.$$



A random state sequence drawn from this MC:

$(2, 2, 2, 1, 1, 2, 2, 1)$

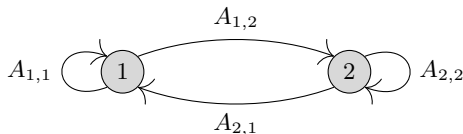
What is the probability of this sequence?

# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\pi = \begin{matrix} & \text{state 1} \\ \text{state 1} & 0.1 \\ \text{state 2} & 0.9 \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & 0.3 & 0.7 \\ \text{state 2} & 0.6 & 0.4 \end{matrix}.$$



A random state sequence drawn from this MC:

(2, 2, 2, 1, 1, 2, 2, 1)

What is the probability of this sequence?

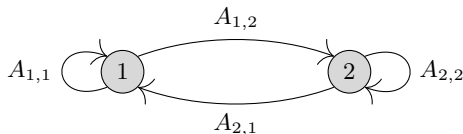
$\pi_2$

# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\pi = \begin{matrix} & \text{state 1} \\ \text{state 1} & 0.1 \\ \text{state 2} & 0.9 \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & 0.3 & 0.7 \\ \text{state 2} & 0.6 & 0.4 \end{matrix}.$$



A random state sequence drawn from this MC:

(2, 2, 2, 1, 1, 2, 2, 1)

What is the probability of this sequence?

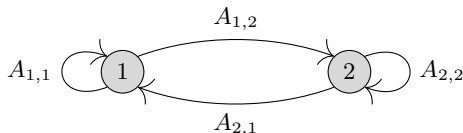
$$\pi_2 \times A_{2,2}$$

# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\pi = \begin{matrix} & \text{state 1} \\ \text{state 1} & 0.1 \\ \text{state 2} & 0.9 \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & 0.3 & 0.7 \\ \text{state 2} & 0.6 & 0.4 \end{matrix}.$$



A random state sequence drawn from this MC:

(2, 2, 2, 1, 1, 2, 2, 1)

What is the probability of this sequence?

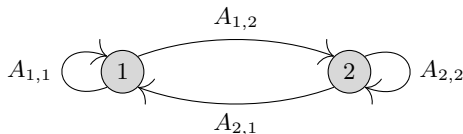
$$\pi_2 \times A_{2,2} \times A_{2,2}$$

# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\pi = \begin{matrix} & \text{state 1} \\ \text{state 1} & \begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix} \\ \text{state 2} & \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix} \\ \text{state 2} & \end{matrix}.$$



A random state sequence drawn from this MC:

(2, 2, **2**, **1**, 1, 2, 2, 1)

What is the probability of this sequence?

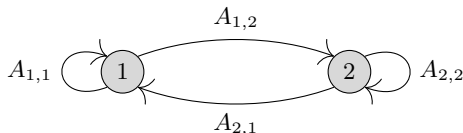
$$\pi_2 \times A_{2,2} \times A_{2,2} \times A_{2,1}$$

# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\pi = \begin{matrix} & \text{state 1} \\ \text{state 1} & 0.1 \\ \text{state 2} & 0.9 \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & 0.3 & 0.7 \\ \text{state 2} & 0.6 & 0.4 \end{matrix}.$$



A random state sequence drawn from this MC:

(2, 2, 2, 1, 1, 2, 2, 1)

What is the probability of this sequence?

$$\pi_2 \times A_{2,2} \times A_{2,2} \times A_{2,1} \times A_{1,1}$$

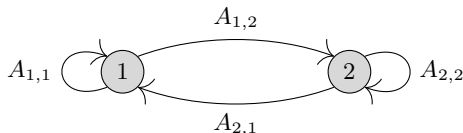


# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\pi = \begin{matrix} & \text{state 1} \\ \text{state 1} & \begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix} \\ \text{state 2} & \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix} \\ \text{state 2} & \end{matrix}.$$



A random state sequence drawn from this MC:

(2, 2, 2, 1, 1, 2, 2, 1)

What is the probability of this sequence?

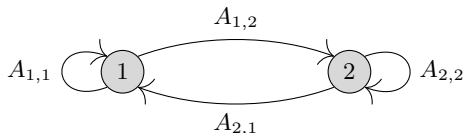
$$\pi_2 \times A_{2,2} \times A_{2,2} \times A_{2,1} \times A_{1,1} \times A_{1,2}$$

# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\pi = \begin{matrix} & \text{state 1} \\ \text{state 1} & 0.1 \\ \text{state 2} & 0.9 \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & 0.3 & 0.7 \\ \text{state 2} & 0.6 & 0.4 \end{matrix}.$$



A random state sequence drawn from this MC:

(2, 2, 2, 1, 1, **2**, **2**, 1)

What is the probability of this sequence?

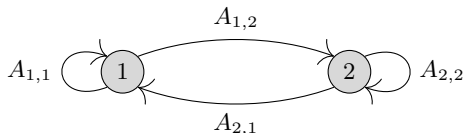
$$\pi_2 \times A_{2,2} \times A_{2,2} \times A_{2,1} \times A_{1,1} \times A_{1,2} \times A_{2,2}$$

# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\pi = \begin{matrix} & \text{state 1} \\ \text{state 1} & 0.1 \\ \text{state 2} & 0.9 \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & 0.3 & 0.7 \\ \text{state 2} & 0.6 & 0.4 \end{matrix}.$$



A random state sequence drawn from this MC:

(2, 2, 2, 1, 1, 2, **2**, **1**)

What is the probability of this sequence?

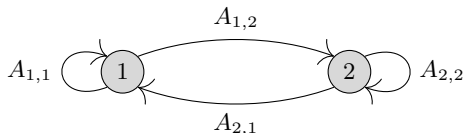
$$\pi_2 \times A_{2,2} \times A_{2,2} \times A_{2,1} \times A_{1,1} \times A_{1,2} \times A_{2,2} \times A_{2,1}$$

# EXAMPLE: A TWO-STATE MARKOV CHAIN

State space:  $\{1, 2\}$ .

Parameters:

$$\pi = \begin{matrix} & \text{state 1} \\ \text{state 1} & \begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix} \\ \text{state 2} & \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \text{state 1} & \text{state 2} \\ \text{state 1} & \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix} \\ \text{state 2} & \end{matrix}.$$



finite machine

A random state sequence drawn from this MC:

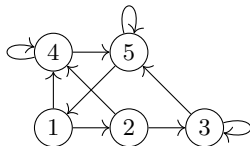
(2, 2, 2, 1, 1, 2, 2, 1)

What is the probability of this sequence?

$$\pi_2 \times A_{2,2} \times A_{2,2} \times A_{2,1} \times A_{1,1} \times A_{1,2} \times A_{2,2} \times A_{2,1} = 0.00435456$$

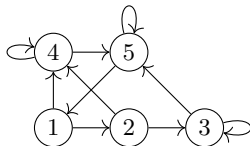
## EXAMPLE: RANDOM WALK ON A DIRECTED GRAPH

Consider a directed graph  $G = (V, E)$  over  $|V| = d$  vertices (self-loops ok).



# EXAMPLE: RANDOM WALK ON A DIRECTED GRAPH

Consider a directed graph  $G = (V, E)$  over  $|V| = d$  vertices (self-loops ok).



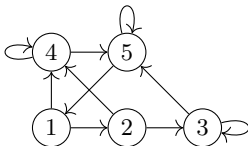
**MC for random walk on  $G$ :**

$$\pi_i = \mathbb{1}\{\text{start vertex is } i\}, \quad A_{i,j} = \frac{\mathbb{1}\{(i,j) \in E\}}{\text{out degree}(i)}.$$

	state 1	state 2	state 3	state 4	state 5
state 1	0	0.5	0	0.5	0
state 2	0	0	0.5	0.5	0
state 3	0	0	0.5	0	0.5
state 4	0	0	0	0.5	0.5
state 5	0.5	0	0	0	0.5

# EXAMPLE: RANDOM WALK ON A DIRECTED GRAPH

Consider a directed graph  $G = (V, E)$  over  $|V| = d$  vertices (self-loops ok).



MC for random walk on  $G$ :

$$\pi_i = \mathbb{1}\{\text{start vertex is } i\}, \quad A_{i,j} = \frac{\mathbb{1}\{(i,j) \in E\}}{\text{out degree}(i)}.$$

? is there are  
weight difference  
on each edge?

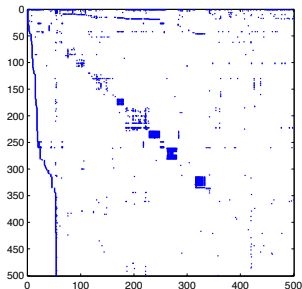
	state 1	state 2	state 3	state 4	state 5
state 1		*		*	
state 2			*	*	
state 3			*		*
state 4				*	*
state 5	*				*

The non-zero pattern of  $A$  gives the adjacency structure of  $G$  (vertices = states).

# EXAMPLE: PAGERANK

**Web graph**  $G = (V, E)$ :

Vertices are webpages, directed edges are hyperlinks between webpages.



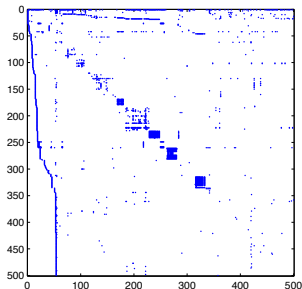
Adjacency matrix of the web graph for 500 web pages.



# EXAMPLE: PAGERANK

**Web graph**  $G = (V, E)$ :

Vertices are webpages, directed edges are hyperlinks between webpages.



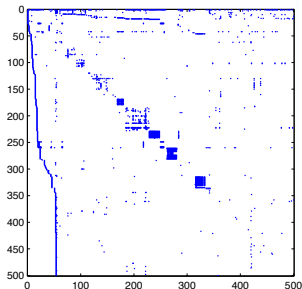
Adjacency matrix of the web graph for 500 web pages.

How popular is webpage  $i$ ?

# EXAMPLE: PAGERANK

**Web graph**  $G = (V, E)$ :

Vertices are webpages, directed edges are hyperlinks between webpages.



Adjacency matrix of the web graph for 500 web pages.

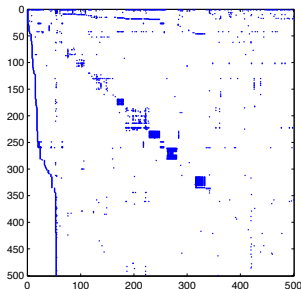
How popular is webpage  $i$ ?

**Possible answer:** probability that **random walk** ends at  $i$  after many steps.

# EXAMPLE: PAGERANK

**Web graph**  $G = (V, E)$ :

Vertices are webpages, directed edges are hyperlinks between webpages.



Adjacency matrix of the web graph for 500 web pages.

How popular is webpage  $i$ ?

**Possible answer:** probability that **random walk** ends at  $i$  after many steps.

$$\Pr(X_t = i) \quad \text{for large } t. \quad \text{may converge}$$

# MARKOV CHAIN STATE DISTRIBUTIONS

# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_2$  in terms of  $\pi$  and  $A$ ?

# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_2$  in terms of  $\pi$  and  $A$ ?

For each  $j \in [d]$ ,

$$\Pr(X_2 = j)$$

# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_2$  in terms of  $\pi$  and  $A$ ?

For each  $j \in [d]$ ,

$$\Pr(X_2 = j) = \sum_{i=1}^d \Pr(X_1 = i, X_2 = j)$$

# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_2$  in terms of  $\pi$  and  $A$ ?

For each  $j \in [d]$ ,

$$\begin{aligned}\Pr(X_2 = j) &= \sum_{i=1}^d \Pr(X_1 = i, X_2 = j) \\ &= \sum_{i=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i)\end{aligned}$$



# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_2$  in terms of  $\pi$  and  $A$ ?

For each  $j \in [d]$ ,

$$\begin{aligned}\Pr(X_2 = j) &= \sum_{i=1}^d \Pr(X_1 = i, X_2 = j) \\ &= \sum_{i=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i) \\ &= \sum_{i=1}^d \pi_i \cdot A_{i,j}\end{aligned}$$

# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_2$  in terms of  $\pi$  and  $A$ ?

For each  $j \in [d]$ ,

$$\begin{aligned}\Pr(X_2 = j) &= \sum_{i=1}^d \Pr(X_1 = i, X_2 = j) \\ &= \sum_{i=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i) \\ &= \sum_{i=1}^d \pi_i \cdot A_{i,j} && \text{Great thinking!!!} \\ &= j\text{-th entry of } \pi^\top A.\end{aligned}$$

# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_3$  in terms of  $\pi$  and  $\mathbf{A}$ ?

# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_3$  in terms of  $\pi$  and  $\mathbf{A}$ ? For each  $k \in [d]$ ,

$$\Pr(X_3 = k) = \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i, X_2 = j, X_3 = k)$$

# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_3$  in terms of  $\pi$  and  $\mathbf{A}$ ? For each  $k \in [d]$ ,

$$\begin{aligned}\Pr(X_3 = k) &= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i, X_2 = j, X_3 = k) \\ &= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i) \cdot \Pr(X_3 = k \mid X_1 = i, X_2 = j)\end{aligned}$$

# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_3$  in terms of  $\pi$  and  $\mathbf{A}$ ? For each  $k \in [d]$ ,

$$\begin{aligned}\Pr(X_3 = k) &= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i, X_2 = j, X_3 = k) \\&= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i) \cdot \Pr(X_3 = k \mid X_1 = i, X_2 = j) \\&= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i) \cdot \Pr(X_3 = k \mid X_2 = j)\end{aligned}$$

# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_3$  in terms of  $\pi$  and  $\mathbf{A}$ ? For each  $k \in [d]$ ,

$$\begin{aligned}\Pr(X_3 = k) &= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i, X_2 = j, X_3 = k) \\&= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i) \cdot \Pr(X_3 = k \mid X_1 = i, X_2 = j) \\&= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i) \cdot \Pr(X_3 = k \mid X_2 = j) \\&= \sum_{i=1}^d \sum_{j=1}^d \pi_i \cdot A_{i,j} \cdot A_{j,k}\end{aligned}$$

# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_3$  in terms of  $\pi$  and  $\mathbf{A}$ ? For each  $k \in [d]$ ,

$$\begin{aligned}\Pr(X_3 = k) &= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i, X_2 = j, X_3 = k) \\&= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i) \cdot \Pr(X_3 = k \mid X_1 = i, X_2 = j) \\&= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i) \cdot \Pr(X_3 = k \mid X_2 = j) \\&= \sum_{i=1}^d \sum_{j=1}^d \pi_i \cdot A_{i,j} \cdot A_{j,k} \\&= \text{\textit{k}-th entry of } \pi^\top \mathbf{A} \mathbf{A}.\end{aligned}$$



# MARGINAL PROBABILITIES

What is the marginal distribution of  $X_3$  in terms of  $\pi$  and  $\mathbf{A}$ ? For each  $k \in [d]$ ,

$$\begin{aligned}\Pr(X_3 = k) &= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i, X_2 = j, X_3 = k) \\&= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i) \cdot \Pr(X_3 = k \mid X_1 = i, X_2 = j) \\&= \sum_{i=1}^d \sum_{j=1}^d \Pr(X_1 = i) \cdot \Pr(X_2 = j \mid X_1 = i) \cdot \Pr(X_3 = k \mid X_2 = j) \\&= \sum_{i=1}^d \sum_{j=1}^d \pi_i \cdot A_{i,j} \cdot A_{j,k} \\&= \text{\textit{k}-th entry of } \pi^\top \mathbf{A} \mathbf{A}.\end{aligned}$$

magic matrix!!!  
It's great!

For any  $t \in \mathbb{N}$ , the marginal distribution of  $X_t$  in terms of  $\pi$  and  $\mathbf{A}$  is

$$\Pr(X_t = k) = \text{\textit{k}-th entry of } \pi^\top \underbrace{\mathbf{A} \mathbf{A} \cdots \mathbf{A}}_{t-1 \text{ times}}$$

# POWERS OF THE TRANSITION MATRIX

The  $(i, j)$ -th entry of  $\mathbf{A}^p = \underbrace{\mathbf{A}\mathbf{A}\cdots\mathbf{A}}_{p \text{ times}}$  is the  $p$ -step transition matrix

$$[\mathbf{A}^p]_{i,j} = \Pr(X_{t+p} = j \mid X_t = i).$$

# POWERS OF THE TRANSITION MATRIX

The  $(i, j)$ -th entry of  $\mathbf{A}^p = \underbrace{\mathbf{A}\mathbf{A}\cdots\mathbf{A}}_{p \text{ times}}$  is the  $p$ -step transition matrix

$$[\mathbf{A}^p]_{i,j} = \Pr(X_{t+p} = j \mid X_t = i).$$

**Example:** State space:  $\{1, 2\}$ . Parameters:  $\boldsymbol{\pi} = \begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix}$ ,  $\mathbf{A} = \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix}$

# POWERS OF THE TRANSITION MATRIX

The  $(i, j)$ -th entry of  $\mathbf{A}^p = \underbrace{\mathbf{A}\mathbf{A}\cdots\mathbf{A}}_{p \text{ times}}$  is the  $p$ -step transition matrix

$$[\mathbf{A}^p]_{i,j} = \Pr(X_{t+p} = j \mid X_t = i).$$

**Example:** State space:  $\{1, 2\}$ . Parameters:  $\boldsymbol{\pi} = \begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix}$ ,  $\mathbf{A} = \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix}$

$$\boldsymbol{\pi}^\top \mathbf{A} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix} = \begin{pmatrix} 0.57 & 0.43 \end{pmatrix}$$

# POWERS OF THE TRANSITION MATRIX

The  $(i, j)$ -th entry of  $\mathbf{A}^p = \underbrace{\mathbf{A}\mathbf{A}\cdots\mathbf{A}}_{p \text{ times}}$  is the  $p$ -step transition matrix

$$[\mathbf{A}^p]_{i,j} = \Pr(X_{t+p} = j \mid X_t = i).$$

**Example:** State space:  $\{1, 2\}$ . Parameters:  $\boldsymbol{\pi} = \begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix}$ ,  $\mathbf{A} = \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix}$

$$\boldsymbol{\pi}^\top \mathbf{A} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix} = \begin{pmatrix} 0.57 & 0.43 \end{pmatrix}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^5 = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.46023 & 0.53977 \\ 0.46266 & 0.53734 \end{pmatrix} = \begin{pmatrix} 0.462417 & 0.537583 \end{pmatrix}$$

# POWERS OF THE TRANSITION MATRIX

The  $(i, j)$ -th entry of  $\mathbf{A}^p = \underbrace{\mathbf{A}\mathbf{A}\cdots\mathbf{A}}_{p \text{ times}}$  is the  $p$ -step transition matrix

$$[\mathbf{A}^p]_{i,j} = \Pr(X_{t+p} = j \mid X_t = i).$$

**Example:** State space:  $\{1, 2\}$ . Parameters:  $\boldsymbol{\pi} = \begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix}$ ,  $\mathbf{A} = \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix}$

$$\boldsymbol{\pi}^\top \mathbf{A} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix} = \begin{pmatrix} 0.57 & 0.43 \end{pmatrix}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^5 = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.46023 & 0.53977 \\ 0.46266 & 0.53734 \end{pmatrix} = \begin{pmatrix} 0.462417 & 0.537583 \end{pmatrix}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^{100} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.461538 & 0.538414 \\ 0.461538 & 0.538414 \end{pmatrix} = \begin{pmatrix} 0.461538 & 0.538414 \end{pmatrix}$$

# POWERS OF THE TRANSITION MATRIX

The  $(i, j)$ -th entry of  $\mathbf{A}^p = \underbrace{\mathbf{A}\mathbf{A}\cdots\mathbf{A}}_{p \text{ times}}$  is the  $p$ -step transition matrix

$$[\mathbf{A}^p]_{i,j} = \Pr(X_{t+p} = j \mid X_t = i).$$

**Example:** State space:  $\{1, 2\}$ . Parameters:  $\boldsymbol{\pi} = \begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix}$ ,  $\mathbf{A} = \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix}$

$$\boldsymbol{\pi}^\top \mathbf{A} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix} = \begin{pmatrix} 0.57 & 0.43 \end{pmatrix}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^5 = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.46023 & 0.53977 \\ 0.46266 & 0.53734 \end{pmatrix} = \begin{pmatrix} 0.462417 & 0.537583 \end{pmatrix}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^{100} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.461538 & 0.538414 \\ 0.461538 & 0.538414 \end{pmatrix} = \begin{pmatrix} 0.461538 & 0.538414 \end{pmatrix}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^{1000} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.461538 & 0.538414 \\ 0.461538 & 0.538414 \end{pmatrix} = \begin{pmatrix} 0.461538 & 0.538414 \end{pmatrix}$$

# POWERS OF THE TRANSITION MATRIX

The  $(i, j)$ -th entry of  $\mathbf{A}^p = \underbrace{\mathbf{A}\mathbf{A}\cdots\mathbf{A}}_{p \text{ times}}$  is the  $p$ -step transition matrix

$$[\mathbf{A}^p]_{i,j} = \Pr(X_{t+p} = j \mid X_t = i).$$

**Example:** State space:  $\{1, 2\}$ . Parameters:  $\boldsymbol{\pi} = \begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix}$ ,  $\mathbf{A} = \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix}$

$$\boldsymbol{\pi}^\top \mathbf{A} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix} = \begin{pmatrix} 0.57 & 0.43 \end{pmatrix}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^5 = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.46023 & 0.53977 \\ 0.46266 & 0.53734 \end{pmatrix} = \begin{pmatrix} 0.462417 & 0.537583 \end{pmatrix}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^{100} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.461538 & 0.538414 \\ 0.461538 & 0.538414 \end{pmatrix} = \begin{pmatrix} 0.461538 & 0.538414 \end{pmatrix}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^{1000} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.461538 & 0.538414 \\ 0.461538 & 0.538414 \end{pmatrix} = \begin{pmatrix} 0.461538 & 0.538414 \end{pmatrix}$$

**Convergence?**



# POWERS OF THE TRANSITION MATRIX

The  $(i, j)$ -th entry of  $\mathbf{A}^p = \underbrace{\mathbf{A}\mathbf{A}\cdots\mathbf{A}}_{p \text{ times}}$  is the  $p$ -step transition matrix

powerful and beautiful!

$$[\mathbf{A}^p]_{i,j} = \Pr(X_{t+p} = j \mid X_t = i).$$

**Example:** State space:  $\{1, 2\}$ . Parameters:  $\boldsymbol{\pi} = \begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix}$ ,  $\mathbf{A} = \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix}$

$$\boldsymbol{\pi}^\top \mathbf{A} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.3 & 0.7 \\ 0.6 & 0.4 \end{pmatrix} = \begin{pmatrix} 0.57 & 0.43 \end{pmatrix} \quad \text{beautiful!}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^5 = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.46023 & 0.53977 \\ 0.46266 & 0.53734 \end{pmatrix} = \begin{pmatrix} 0.462417 & 0.537583 \end{pmatrix}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^{100} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.461538 & 0.538414 \\ 0.461538 & 0.538414 \end{pmatrix} = \begin{pmatrix} 0.461538 & 0.538414 \end{pmatrix}$$

$$\boldsymbol{\pi}^\top \mathbf{A}^{1000} = \begin{pmatrix} 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.461538 & 0.538414 \\ 0.461538 & 0.538414 \end{pmatrix} = \begin{pmatrix} 0.461538 & 0.538414 \end{pmatrix}$$

Convergence? Doesn't even seem to matter what  $\boldsymbol{\pi}$  is!

# LIMITING STATE DISTRIBUTION

Certain “nice” transition matrices  $\mathbf{A} \in \mathbb{R}^{d \times d}$  have the property that

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \text{stochastic matrix with identical rows}$$

# LIMITING STATE DISTRIBUTION

Certain “nice” transition matrices  $\mathbf{A} \in \mathbb{R}^{d \times d}$  have the property that

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \text{stochastic matrix with identical rows} =: \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix}$$

# LIMITING STATE DISTRIBUTION

Certain “nice” transition matrices  $\mathbf{A} \in \mathbb{R}^{d \times d}$  have the property that

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \text{stochastic matrix with identical rows} =: \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix}$$

What can we say about  $\mathbf{q}$ ?

# LIMITING STATE DISTRIBUTION

Certain “nice” transition matrices  $\mathbf{A} \in \mathbb{R}^{d \times d}$  have the property that

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \text{stochastic matrix with identical rows} =: \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix}$$

What can we say about  $\mathbf{q}$ ? For such  $\mathbf{A}$ ,

$$\lim_{p \rightarrow \infty} \mathbf{A}^p$$

# LIMITING STATE DISTRIBUTION

Certain “nice” transition matrices  $\mathbf{A} \in \mathbb{R}^{d \times d}$  have the property that

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \text{stochastic matrix with identical rows} =: \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix}$$

What can we say about  $\mathbf{q}$ ? For such  $\mathbf{A}$ ,

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \left( \lim_{p \rightarrow \infty} \mathbf{A}^{p-1} \right) \mathbf{A}$$

# LIMITING STATE DISTRIBUTION

Certain “nice” transition matrices  $\mathbf{A} \in \mathbb{R}^{d \times d}$  have the property that

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \text{stochastic matrix with identical rows} =: \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix}$$

What can we say about  $\mathbf{q}$ ? For such  $\mathbf{A}$ ,

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \left( \lim_{p \rightarrow \infty} \mathbf{A}^{p-1} \right) \mathbf{A} = \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix} \mathbf{A}$$

# LIMITING STATE DISTRIBUTION

Certain “nice” transition matrices  $\mathbf{A} \in \mathbb{R}^{d \times d}$  have the property that

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \text{stochastic matrix with identical rows} =: \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix}$$

What can we say about  $\mathbf{q}$ ? For such  $\mathbf{A}$ ,

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \left( \lim_{p \rightarrow \infty} \mathbf{A}^{p-1} \right) \mathbf{A} = \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix} \mathbf{A} = \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix}$$



# LIMITING STATE DISTRIBUTION

Certain “nice” transition matrices  $\mathbf{A} \in \mathbb{R}^{d \times d}$  have the property that

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \text{stochastic matrix with identical rows} =: \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix}$$

What can we say about  $\mathbf{q}$ ? For such  $\mathbf{A}$ ,

$$\lim_{p \rightarrow \infty} \mathbf{A}^p = \left( \lim_{p \rightarrow \infty} \mathbf{A}^{p-1} \right) \mathbf{A} = \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix} \mathbf{A} = \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix}$$

i.e.,

$$\mathbf{q}^\top \mathbf{A} = \mathbf{q}^\top. \quad (\star)$$

# LIMITING STATE DISTRIBUTION

Certain “nice” transition matrices  $A \in \mathbb{R}^{d \times d}$  have the property that

$\lim_{p \rightarrow \infty} A^p =$  stochastic matrix with identical rows  $=:$

$$\begin{pmatrix} \text{---} & q^\top & \text{---} \\ \text{---} & q^\top & \text{---} \\ & \vdots & \\ \text{---} & q^\top & \text{---} \end{pmatrix}$$

What can we say about  $q$ ? For such  $A$ ,

$$\lim_{p \rightarrow \infty} A^p = \left( \lim_{p \rightarrow \infty} A^{p-1} \right) A = \begin{pmatrix} \text{---} & q^\top & \text{---} \\ \text{---} & q^\top & \text{---} \\ & \vdots & \\ \text{---} & q^\top & \text{---} \end{pmatrix} A = \begin{pmatrix} \text{---} & q^\top & \text{---} \\ \text{---} & q^\top & \text{---} \\ & \vdots & \\ \text{---} & q^\top & \text{---} \end{pmatrix}$$

try to find out the  $q$  !!!

i.e.,

$$q^\top A = q^\top. \quad (\star)$$

A solution  $q$  to  $(\star)$ , is called a **stationary distribution**.

# STATIONARY DISTRIBUTION

Suppose a MC has a unique stationary distribution  $\mathbf{q} = (q_1, q_2, \dots, q_d)$ .

# STATIONARY DISTRIBUTION

Suppose a MC has a unique stationary distribution  $\mathbf{q} = (q_1, q_2, \dots, q_d)$ .

- For any  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \Pr \left( \left| \frac{1}{n} \sum_{t=1}^n \mathbb{1}\{X_t = i\} - q_i \right| > \varepsilon \right) \rightarrow 0.$$

Law of Large Numbers for MCs.

# STATIONARY DISTRIBUTION

Suppose a MC has a unique stationary distribution  $\mathbf{q} = (q_1, q_2, \dots, q_d)$ .

- For any  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \Pr \left( \left| \frac{1}{n} \sum_{t=1}^n \mathbb{1}\{X_t = i\} - q_i \right| > \varepsilon \right) \rightarrow 0.$$

Law of Large Numbers for MCs.

- However, rate of convergence **not the same as in the iid case**.

# STATIONARY DISTRIBUTION

Suppose a MC has a unique stationary distribution  $\mathbf{q} = (q_1, q_2, \dots, q_d)$ .

- For any  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \Pr \left( \left| \frac{1}{n} \sum_{t=1}^n \mathbb{1}\{X_t = i\} - q_i \right| > \varepsilon \right) \rightarrow 0.$$

Law of Large Numbers for MCs.

- However, rate of convergence **not the same as in the iid case**.

Critically depends on how quickly  $\Pr(X_t = \cdot) \rightarrow \mathbf{q}$  (**mixing rate**).

# STATIONARY DISTRIBUTION

Suppose a MC has a unique stationary distribution  $\mathbf{q} = (q_1, q_2, \dots, q_d)$ .

- For any  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \Pr \left( \left| \frac{1}{n} \sum_{t=1}^n \mathbb{1}\{X_t = i\} - q_i \right| > \varepsilon \right) \rightarrow 0.$$

as the size of n growing  
up, the empirical  
appearance is getting  
close to  $q_i$

Law of Large Numbers for MCs.

- However, rate of convergence **not the same as in the iid case.**

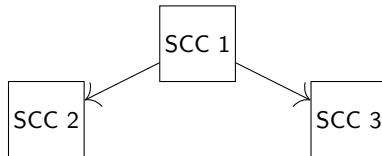
not like chernoff  
bound

Critically depends on how quickly  $\Pr(X_t = \cdot) \rightarrow \mathbf{q}$  (**mixing rate**).

**When does a MC even have a unique stationary distribution?**

# WHAT CAN GO WRONG

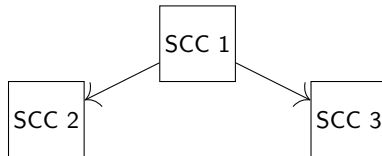
1. Directed graph underlying  $A$  has more than one strongly connected component “sinks”  $\rightarrow$  stationary distribution may not be unique.





# WHAT CAN GO WRONG

1. Directed graph underlying  $A$  has more than one strongly connected component “sinks”  $\rightarrow$  stationary distribution may not be unique.



Markov chains with only one strongly connected component are called **irreducible**.

# WHAT CAN GO WRONG

2. Oscillation among two or more states  $\rightarrow$  limit does not exist.

# WHAT CAN GO WRONG

2. **Oscillation among two or more states**  $\rightarrow$  **limit does not exist.**

Example:



If start at state 1, then never at state 1 on even time steps.

# WHAT CAN GO WRONG

2. **Oscillation among two or more states**  $\rightarrow$  **limit does not exist.**

Example:



If start at state 1, then never at state 1 on even time steps.

Markov chains without such oscillation are called **aperiodic**.

(Formally: there exists  $p_0$  s.t. for all  $p \geq p_0$ ,  $[A^p]_{i,i} > 0$  for all  $i \in [d]$ .)

# WHAT CAN GO WRONG

## 2. Oscillation among two or more states $\rightarrow$ limit does not exist.

Example:



If start at state 1, then never at state 1 on even time steps.

Markov chains without such oscillation are called **aperiodic**.

(Formally: there exists  $p_0$  s.t. for all  $p \geq p_0$ ,  $[A^p]_{i,i} > 0$  for all  $i \in [d]$ .)

If every state  $i \in [d]$  has  $A_{i,i} > 0$ , then aperiodicity is guaranteed.

# CONDITIONS FOR UNIQUE STATIONARY DISTRIBUTION

**Theorem:** If MC with transition matrix  $A$  is *irreducible* and *aperiodic*, then

- ▶ There is a unique stationary distribution  $\mathbf{q}$  (which satisfies  $\mathbf{q}^\top A = \mathbf{q}^\top$ ).

- ▶  $\lim_{p \rightarrow \infty} A^p = \begin{pmatrix} \text{---} & \mathbf{q}^\top & \text{---} \\ \text{---} & \mathbf{q}^\top & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{q}^\top & \text{---} \end{pmatrix}.$

# COMPUTING THE STATIONARY DISTRIBUTION

For irreducible and aperiodic MCs, the  $\mathbf{q}$  that satisfies

$$\mathbf{q}^\top \mathbf{A} = \mathbf{q}^\top$$

is unique.

# COMPUTING THE STATIONARY DISTRIBUTION

For irreducible and aperiodic MCs, the  $\mathbf{q}$  that satisfies

$$\mathbf{q}^\top \mathbf{A} = \mathbf{q}^\top$$

is unique. Therefore, suffices to find *left eigenvector* of  $\mathbf{A}$  with eigenvalue 1.



# COMPUTING THE STATIONARY DISTRIBUTION

For irreducible and aperiodic MCs, the  $\mathbf{q}$  that satisfies

$$\mathbf{q}^\top \mathbf{A} = \mathbf{q}^\top$$

is unique. Therefore, suffices to find *left eigenvector* of  $\mathbf{A}$  with eigenvalue 1.  
In fact,  $\mathbf{A}$  has no other eigenvalue of larger modulus!

# COMPUTING THE STATIONARY DISTRIBUTION

For irreducible and aperiodic MCs, the  $\mathbf{q}$  that satisfies

$$\mathbf{q}^\top \mathbf{A} = \mathbf{q}^\top$$

is unique. Therefore, suffices to find *left eigenvector* of  $\mathbf{A}$  with eigenvalue 1.  
In fact,  $\mathbf{A}$  has no other eigenvalue of larger modulus!

**Direct method:** Find any vector in *left null space* of  $\mathbf{A} - \mathbf{I}$

$$\mathbf{q}^\top (\mathbf{A} - \mathbf{I}) = \mathbf{0},$$

and properly normalize it to be a state distribution.

# COMPUTING THE STATIONARY DISTRIBUTION

For irreducible and aperiodic MCs, the  $q$  that satisfies

$$q^\top A = q^\top$$

is unique. Therefore, suffices to find *left eigenvector of  $A$  with eigenvalue 1.*

In fact,  $A$  has no other eigenvalue of larger modulus!

**Direct method:** Find any vector in *left null space* of  $A - I$

$$q^\top (A - I) = 0,$$

and properly normalize it to be a state distribution.

**Power method:**

**initialize**  $q$  arbitrarily.

**repeat**

$q^\top := q^\top A.$  A is usually sparse

**until** bored.

**return**  $q.$

# EXAMPLE: PAGERANK

## Random walk on web graph:

- ▶ definitely **not irreducible**,  
(some pages have no links to other pages);
- ▶ probably **not aperiodic**.

# EXAMPLE: PAGERANK

**Random walk on web graph:**

- ▶ definitely **not irreducible**,  
(some pages have no links to other pages);
- ▶ probably **not aperiodic**.

**Modification:**

$$\tilde{\mathbf{A}} := (1 - \alpha)\mathbf{A} + \frac{\alpha}{d} \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{pmatrix}.$$

New MC (with  $\tilde{\mathbf{A}}$ ) is both irreducible and aperiodic.

# EXAMPLE: PAGERANK

## Random walk on web graph:

- ▶ definitely **not irreducible.**  
(some pages have no links to other pages);
- ▶ probably **not aperiodic.**

too many sinks~

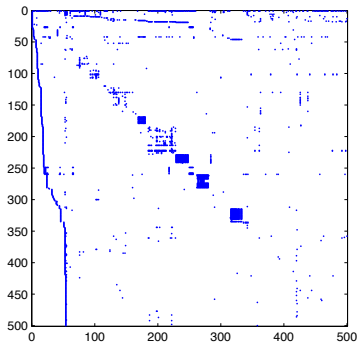
## Modification:

$$\tilde{\mathbf{A}} := (1 - \alpha)\mathbf{A} + \frac{\alpha}{d} \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{pmatrix}.$$

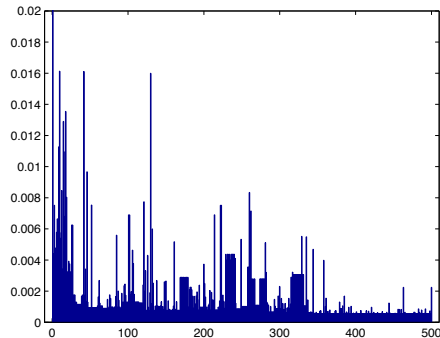
New MC (with  $\tilde{\mathbf{A}}$ ) is both irreducible and aperiodic.

PageRank scores = stationary distribution of this new MC.

# EXAMPLE: PAGERANK



Adjacency matrix of the web graph  
for 500 web pages.



PageRank distribution.

(From K. Murphy, "Machine Learning", MIT Press 2012.)

## EXAMPLE: SEMI-SUPERVISED LEARNING

Have some **labeled data**  $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$  from  $\mathcal{X} \times \{\pm 1\}$ , and also many **unlabeled data**  $x_{m+1}, x_{m+1}, \dots, x_{m+n}$  from  $\mathcal{X}$ .



## EXAMPLE: SEMI-SUPERVISED LEARNING

Have some **labeled data**  $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$  from  $\mathcal{X} \times \{\pm 1\}$ , and also many **unlabeled data**  $x_{m+1}, x_{m+1}, \dots, x_{m+n}$  from  $\mathcal{X}$ .

**Main question:** How to take advantage of  $U$ ?

## EXAMPLE: SEMI-SUPERVISED LEARNING

Have some **labeled data**  $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$  from  $\mathcal{X} \times \{\pm 1\}$ , and also many **unlabeled data**  $x_{m+1}, x_{m+1}, \dots, x_{m+n}$  from  $\mathcal{X}$ .

**Main question:** How to take advantage of  $U$ ?

[Zhu, Ghahramani, and Lafferty, 2003]

- ▶ Construct weighted similarity graph  $G = (V, W)$  over all data.

# EXAMPLE: SEMI-SUPERVISED LEARNING

Have some **labeled data**  $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$  from  $\mathcal{X} \times \{\pm 1\}$ , and also many **unlabeled data**  $x_{m+1}, x_{m+1}, \dots, x_{m+n}$  from  $\mathcal{X}$ .

**Main question:** How to take advantage of  $U$ ?

[Zhu, Ghahramani, and Lafferty, 2003]

- ▶ Construct weighted similarity graph  $G = (V, W)$  over all data.

For example:

- ▶  $V = \{1, 2, \dots, m + n\}$ .
- ▶ Weight  $W_{i,j} = \exp\left(-\frac{1}{2} \text{dist}(x_i, x_j)^2\right)$ .

# EXAMPLE: SEMI-SUPERVISED LEARNING

Have some **labeled data**  $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$  from  $\mathcal{X} \times \{\pm 1\}$ , and also many **unlabeled data**  $x_{m+1}, x_{m+1}, \dots, x_{m+n}$  from  $\mathcal{X}$ .

**Main question:** How to take advantage of  $U$ ?

[Zhu, Ghahramani, and Lafferty, 2003]

- ▶ Construct weighted similarity graph  $G = (V, W)$  over all data.

For example:

- ▶  $V = \{1, 2, \dots, m+n\}$ .
- ▶ Weight  $W_{i,j} = \exp\left(-\frac{1}{2} \text{dist}(x_i, x_j)^2\right)$ .
- ▶ Weighted random walk MC:

$$A_{i,j} = \frac{W_{i,j}}{\sum_{k=1}^{m+n} W_{i,k}}.$$

# EXAMPLE: SEMI-SUPERVISED LEARNING

Have some **labeled data**  $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$  from  $\mathcal{X} \times \{\pm 1\}$ , and also many **unlabeled data**  $x_{m+1}, x_{m+1}, \dots, x_{m+n}$  from  $\mathcal{X}$ .

**Main question:** How to take advantage of  $U$ ?

$U$  : unlabeled data

[Zhu, Ghahramani, and Lafferty, 2003]

- ▶ Construct weighted similarity graph  $G = (V, W)$  over all data.

For example:

- ▶  $V = \{1, 2, \dots, m+n\}$ .
- ▶ Weight  $W_{i,j} = \exp\left(-\frac{1}{2} \text{dist}(x_i, x_j)^2\right)$ .
- ▶ Weighted random walk MC:

$$A_{i,j} = \frac{W_{i,j}}{\sum_{k=1}^{m+n} W_{i,k}}.$$

normalize (equal to 0)

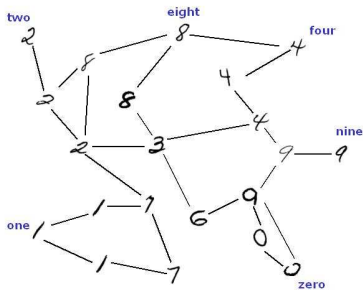
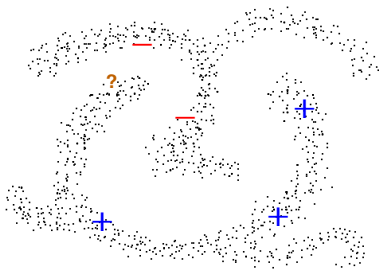
the large the distance, the smaller the weight

- ▶ Start weighted random walk starting from **unlabeled point**  $x_{m+i}$ .

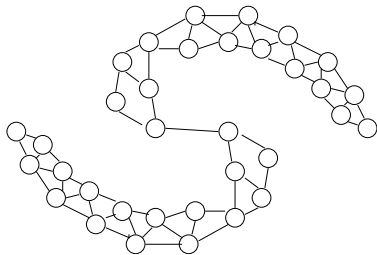
If **first labeled point reached** has label  $y \in \{\pm 1\}$ , then use  $\hat{y}_{m+i} := y$  as the label for  $x_{m+i}$ .

(Can actually compute, in closed form, the probabilities of  $\hat{y}_{m+i} = y$  for each  $y$ .)

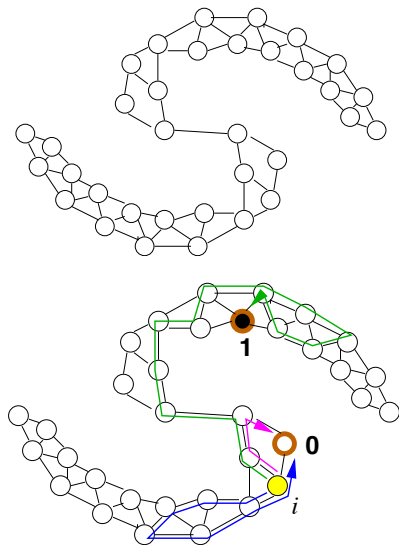
# EXAMPLE: SEMI-SUPERVISED LEARNING



# EXAMPLE: SEMI-SUPERVISED LEARNING



# EXAMPLE: SEMI-SUPERVISED LEARNING





- ▶ Markov property: past and future are conditionally independent given the present.
- ▶ Transition matrix: the conditional next-state distributions for each state.
- ▶ Random walk on graphs: extremely important process, very well-studied, many applications (including in ML, statistics, etc).
- ▶ **Irreducible and aperiodic Markov chains have limiting behavior:**  
doesn't matter where you start, eventually marginal state distribution is the stationary distribution.

Some qualities similar to iid processes, some rather different.

Related to eigenvectors/eigenvalues, computation via power method.

- ▶ Forms the basis of PageRank.