# 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

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

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

# 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 $\{X_\tau : \tau \le t\}$ only depends on the value of the current state $X_t$.

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

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

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

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

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

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

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 $[d] := \{1, 2, \ldots, d\}$.

# 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 $[d] := \{1, 2, \ldots, d\}$.

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

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

# 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 $[d] := \{1, 2, \ldots, d\}$.

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

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

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

$$A_{i,j} \;=\; \Pr(X_{t+1} = j \mid X_t = i)$$
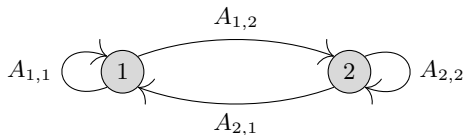
  (rows of $\boldsymbol{A}$ are probability vectors).

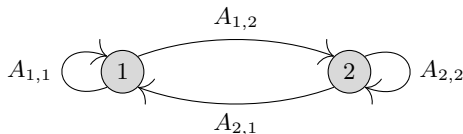  Also called a *transition matrix* or *(right) stochastic matrix*.

# Example: a two-state Markov chain

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

$$\boldsymbol{\pi} = \begin{array}{c} \text{state 1} \\ \text{state 2} \end{array}\!\!\begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix}, \quad \boldsymbol{A} = \begin{array}{c} \\ \text{state 1} \\ \text{state 2} \end{array}\!\!\begin{pmatrix} \overset{\text{state 1}}{0.3} & \overset{\text{state 2}}{0.7} \\ 0.6 & 0.4 \end{pmatrix}.$$

# Example: a two-state Markov chain

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

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



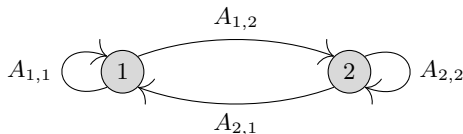A random state sequence drawn from this MC:

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

# EXAMPLE: A TWO-STATE MARKOV CHAIN

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

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



A random state sequence drawn from this MC:
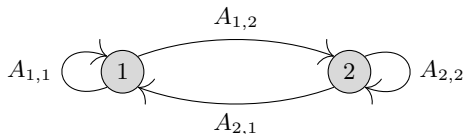
$$(2, 2, 2, 1, 1, 2, 2, 1)$$

What is the probability of this sequence?

# EXAMPLE: A TWO-STATE MARKOV CHAIN

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

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



A random state sequence drawn from this MC:

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

What is the probability of this sequence?

$\pi_2$

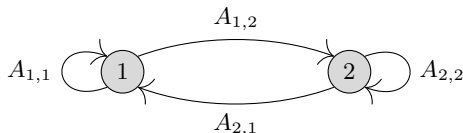# Example: a two-state Markov chain

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

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



A random state sequence drawn from this MC:

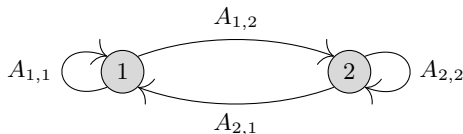$$(2, 2, 2, 1, 1, 2, 2, 1)$$

What is the probability of this sequence?

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

# EXAMPLE: A TWO-STATE MARKOV CHAIN

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

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



A random state sequence drawn from this MC:

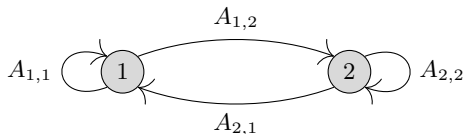$$(2, 2, 2, 1, 1, 2, 2, 1)$$

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:

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



A random state sequence drawn from this MC:
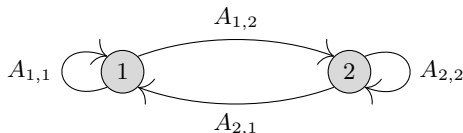
$$(2, 2, 2, 1, 1, 2, 2, 1)$$

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:

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



A random state sequence drawn from this MC:

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

What is the probability of this sequence?

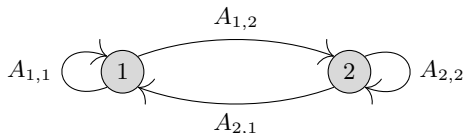$\pi_2 \times A_{2,2} \times A_{2,2} \times A_{2,1} \times A_{1,1}$

# EXAMPLE: A TWO-STATE MARKOV CHAIN

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

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



A random state sequence drawn from this MC:

$$(2, 2, 2, 1, 1, 2, 2, 1)$$
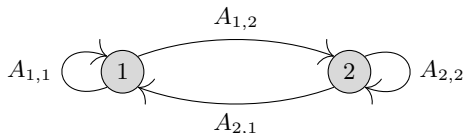
What is the probability of this sequence?

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

# EXAMPLE: A TWO-STATE MARKOV CHAIN

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

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



A random state sequence drawn from this MC:

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

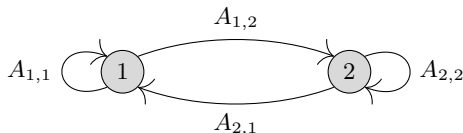What is the probability of this sequence?

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

# EXAMPLE: A TWO-STATE MARKOV CHAIN

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

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



A random state sequence drawn from this MC:

$$(2, 2, 2, 1, 1, 2, 2, 1)$$
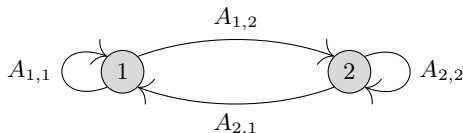
What is the probability of this sequence?

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

# Example: a two-state Markov chain

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

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



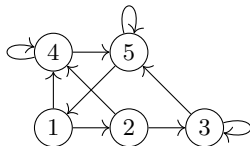A random state sequence drawn from this MC:

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

What is the probability of this sequence?

$$\pi_2 \times A_{2,2} \times A_{2,2} \times A_{2,1} \times A_{1,1} \times A_{1,2} \times A_{2,2} \times A_{2,1} = 0.00435456$$
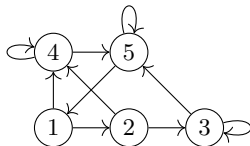
# EXAMPLE: RANDOM WALK ON A DIRECTED GRAPH

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 $G$:**

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

$$
\begin{array}{c}
\\
\text{state 1} \\
\text{state 2} \\
\text{state 3} \\
\text{state 4} \\
\text{state 5}
\end{array}
\begin{array}{ccccc}
\text{state 1} & \text{state 2} & \text{state 3} & \text{state 4} & \text{state 5} \\
\left( \begin{array}{ccccc}
0 & 0.5 & 0 & 0.5 & 0 \\
0 & 0 & 0.5 & 0.5 & 0 \\
0 & 0 & 0.5 & 0 & 0.5 \\
0 & 0 & 0 & 0.5 & 0.5 \\
0.5 & 0 & 0 & 0 & 0.5
\end{array} \right)
\end{array}
$$

# EXAMPLE: RANDOM WALK ON A DIRECTED GRAPH

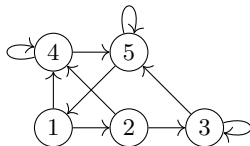Consider a directed graph $G = (V, E)$ over $|V| = d$ vertices (self-loops ok).



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

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

$$
\begin{array}{c}
\\
\text{state 1} \\
\text{state 2} \\
\text{state 3} \\
\text{state 4} \\
\text{state 5}
\end{array}
\begin{array}{ccccc}
\text{state 1} & \text{state 2} & \text{state 3} & \text{state 4} & \text{state 5} \\
\left(\begin{array}{ccccc}
 & * & & * & \\
 & & * & * & \\
 & & * & & * \\
 & & & * & * \\
* & & & & *
\end{array}\right)
\end{array}
$$

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.



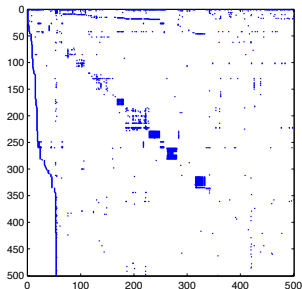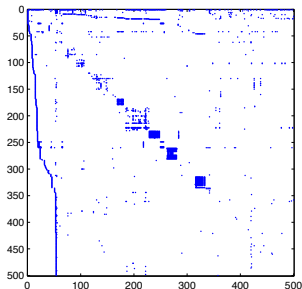Adjacency matrix of the web graph for 500 web pages.

How popular is webpage $i$?

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

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

# Markov chain state distributions

# Marginal probabilities

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

# Marginal probabilities

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

For each $j \in [d]$,

$$\Pr(X_2 = j)$$

# Marginal probabilities

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

For each $j \in [d]$,

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

# Marginal probabilities

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

For each $j \in [d]$,

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

# Marginal probabilities

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

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 $\boldsymbol{\pi}$ and $\boldsymbol{A}$?

For each $j \in [d]$,

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

# Marginal probabilities

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

## Marginal probabilities

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

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

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

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

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

# Marginal probabilities

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

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

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

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

# Marginal probabilities

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

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

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

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

$$= \sum_{i=1}^{d} \sum_{j=1}^{d} \pi_i \cdot A_{i,j} \cdot A_{j,k}$$

# Marginal probabilities

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

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

# Marginal probabilities

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

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

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

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

# POWERS OF THE TRANSITION MATRIX

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

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

# POWERS OF THE TRANSITION MATRIX

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

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

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

# POWERS OF THE TRANSITION MATRIX

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

$$\left[\boldsymbol{A}^p\right]_{i,j} = \Pr(X_{t+p} = j \,|\, X_t = i).$$

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

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

# POWERS OF THE TRANSITION MATRIX

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

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

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

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

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

# POWERS OF THE TRANSITION MATRIX

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

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

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

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

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

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

# POWERS OF THE TRANSITION MATRIX

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

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

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

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

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

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

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

# POWERS OF THE TRANSITION MATRIX

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

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

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

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

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

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

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

**Convergence?**

# POWERS OF THE TRANSITION MATRIX

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

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

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

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

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

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

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

**Convergence? Doesn't even seem to matter what $\pi$ is!**

# Limiting state distribution

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

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

# Limiting state distribution

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

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

# LIMITING STATE DISTRIBUTION

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

$$\lim_{p \to \infty} \boldsymbol{A}^p = \text{stochastic matrix with identical rows} =: \begin{pmatrix} \underline{\quad} & \boldsymbol{q}^\top & \underline{\quad} \\ \underline{\quad} & \boldsymbol{q}^\top & \underline{\quad} \\ & \vdots & \\ \underline{\quad} & \boldsymbol{q}^\top & \underline{\quad} \end{pmatrix}$$

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

# LIMITING STATE DISTRIBUTION

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

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

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

$$\lim_{p \to \infty} \boldsymbol{A}^p$$

# LIMITING STATE DISTRIBUTION

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

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

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

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

# LIMITING STATE DISTRIBUTION

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

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

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

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

# LIMITING STATE DISTRIBUTION

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

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

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

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

# LIMITING STATE DISTRIBUTION

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

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

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

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

i.e.,

$$\boldsymbol{q}^\top \boldsymbol{A} = \boldsymbol{q}^\top. \tag{$\star$}$$

# LIMITING STATE DISTRIBUTION

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

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

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

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

i.e.,

$$\boldsymbol{q}^\top \boldsymbol{A} \; = \; \boldsymbol{q}^\top. \tag{$\star$}$$

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

# STATIONARY DISTRIBUTION

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

# STATIONARY DISTRIBUTION

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

- For any $\varepsilon > 0$,

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

Law of Large Numbers for MCs.

# STATIONARY DISTRIBUTION

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

- For any $\varepsilon > 0$,

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

Law of Large Numbers for MCs.

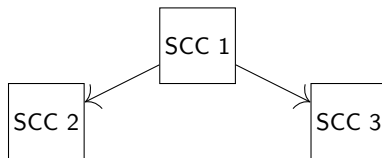- However, rate of convergence not the same as in the iid case.

# STATIONARY DISTRIBUTION

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

- For any $\varepsilon > 0$,

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

  Law of Large Numbers for MCs.

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

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

# STATIONARY DISTRIBUTION

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

- For any $\varepsilon > 0$,

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

  Law of Large Numbers for MCs.

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

  Critically depends on how quickly $\Pr(X_t = \cdot) \to \boldsymbol{q}$ (**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"** $\longrightarrow$ stationary distribution may not be unique.

# WHAT CAN GO WRONG

1. **Directed graph underlying $A$ has more than one strongly connected component "sinks"** $\longrightarrow$ 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** $\longrightarrow$ limit does not exist.

# What can go wrong

2. **Oscillation among two or more states** $\longrightarrow$ 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** $\longrightarrow$ 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 $p \geq p_0$, $\left[\boldsymbol{A}^p\right]_{i,i} > 0$ for all $i \in [d]$.)

# WHAT CAN GO WRONG

2. **Oscillation among two or more states** $\longrightarrow$ 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 $p \geq p_0$, $\left[\boldsymbol{A}^p\right]_{i,i} > 0$ 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 $q$ (which satisfies $q^\top A = q^\top$).

- $\lim_{p \to \infty} A^p = \begin{pmatrix} \rule[0.5ex]{1em}{0.4pt} & q^\top & \rule[0.5ex]{1em}{0.4pt} \\ \rule[0.5ex]{1em}{0.4pt} & q^\top & \rule[0.5ex]{1em}{0.4pt} \\ & \vdots & \\ \rule[0.5ex]{1em}{0.4pt} & q^\top & \rule[0.5ex]{1em}{0.4pt} \end{pmatrix}.$

# Computing the stationary distribution

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

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

is unique.

# 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 $1$.

# 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!

# 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.

# 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**
>     $q^\top := q^\top A$.
> **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**.
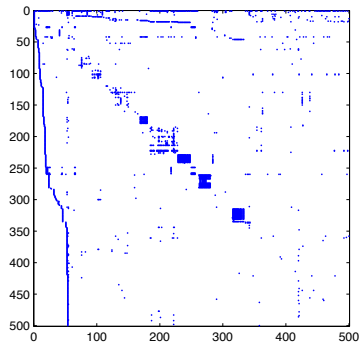
**Modification**:

$$\widetilde{\boldsymbol{A}} \;:=\; (1-\alpha)\boldsymbol{A} \;+\; \frac{\alpha}{d} \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{pmatrix}.$$
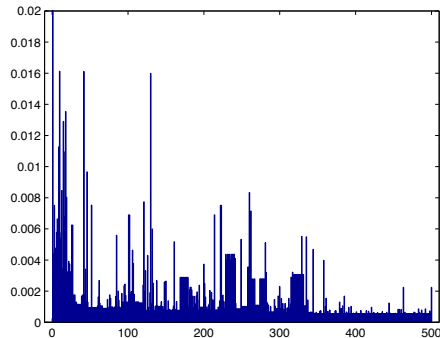
New MC (with $\widetilde{\boldsymbol{A}}$) is both irreducible and aperiodic.

# EXAMPLE: PAGERANK

**Random walk on web graph**:

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

**Modification**:

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

New MC (with $\widetilde{\boldsymbol{A}}$) 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** $(x_1, y_1), (x_2, y_2), \ldots, (x_m, y_m)$ from $\mathcal{X} \times \{\pm 1\}$, and also many **unlabeled data** $x_{m+1}, x_{m+1}, \ldots, x_{m+n}$ from $\mathcal{X}$.

# EXAMPLE: SEMI-SUPERVISED LEARNING

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

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

# EXAMPLE: SEMI-SUPERVISED LEARNING

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

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

# EXAMPLE: SEMI-SUPERVISED LEARNING

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

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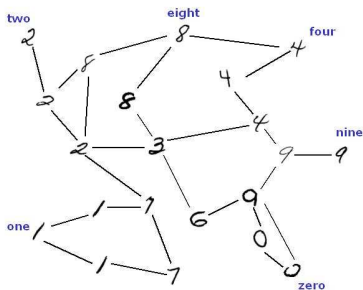
▶ Weighted random walk MC:

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

# EXAMPLE: SEMI-SUPERVISED LEARNING

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

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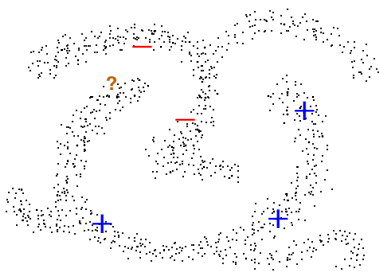
▶ Weighted random walk MC:

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

▶ Start weighted random walk starting from **unlabeled point** $x_{m+i}$.
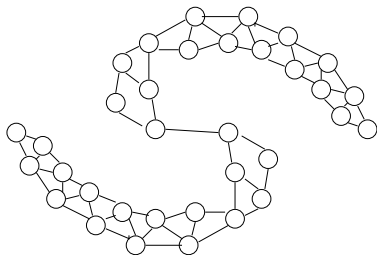
If **first labeled point reached** has label $y \in \{\pm 1\}$, then use $\hat{y}_{m+i} := y$ as the label for $x_{m+i}$.

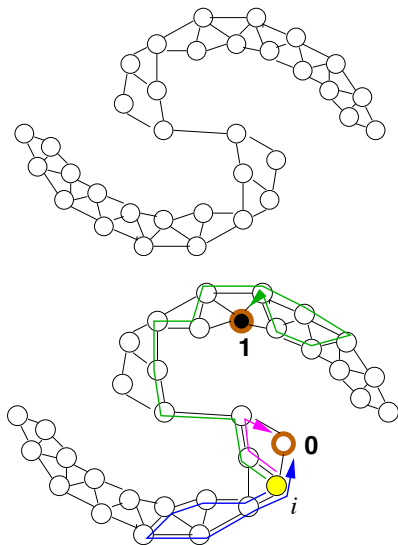(Can actually compute, in closed form, the probabilities of $\hat{y}_{m+i} = y$ for each $y$.)

# EXAMPLE: SEMI-SUPERVISED LEARNING

# Example: semi-supervised learning

# Example: semi-supervised learning

# Recap

- 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.