

# CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

April 2, 2021

1 / 19

Logistics

## Outline

1 Logistics

2 (Hidden) Markov models I

3 / 19

## Outline

1 Logistics

2 (Hidden) Markov models I

2 / 19

Logistics

## Logistics

- **April 7, 2021** is a Wellness day, there will be no class.

4 / 19

## Outline

### 1 Logistics

### 2 (Hidden) Markov models I

- Markov chain
- Hidden Markov Model

5 / 19

## Markov Models

Markov models are powerful probabilistic tools to analyze **sequential data**:

- text or speech data
- stock market data
- gene data
- ...

6 / 19

## Definition

A **Markov chain** is a stochastic process with **Markov property**: a sequence of random variables  $Z_1, Z_2, \dots$  s.t.

$$P(Z_{t+1} \mid Z_{1:t}) = P(Z_{t+1} \mid Z_t) \quad (\text{Markov property})$$

i.e. *the current state only depends on the most recent state* (notation  $Z_{1:t}$  denotes the sequence  $Z_1, \dots, Z_t$ ).

We only consider the following case:

- All  $Z_t$ 's take value from the same **discrete** set  $\{1, \dots, S\}$
- $P(Z_{t+1} = s' \mid Z_t = s) = a_{s,s'}$ , known as **transition probability**
- $P(Z_1 = s) = \pi_s$
- $(\{\pi_s\}, \{a_{s,s'}\}) = (\boldsymbol{\pi}, \mathbf{A})$  are **parameters of the model**

7 / 19

## Examples

- Example 1 (**Language model**)

States  $[S]$  represent a dictionary of words,

$$a_{\text{ice,cream}} = P(Z_{t+1} = \text{cream} \mid Z_t = \text{ice})$$

is an example of the transition probability.

- Example 2 (**Weather**)

States  $[S]$  represent weather at each day

$$a_{\text{sunny,rainy}} = P(Z_{t+1} = \text{rainy} \mid Z_t = \text{sunny})$$

8 / 19

## High-order Markov chain

*Is the Markov assumption reasonable?* Not completely for the language model for example.

Higher order Markov chains make it more reasonable, e.g.

$$P(Z_{t+1} | Z_{1:t}) = P(Z_{t+1} | Z_t, Z_{t-1}) \quad (\text{second-order Markov})$$

i.e. the current word only depends on the last two words.

Learning higher order Markov chains is similar, but more expensive.

We only consider standard Markov chains.

## Learning from examples

Now suppose we have observed  $N$  sequences of examples, say  $\mathcal{D}$ :

- $z_{1,1}, \dots, z_{1,T}$
- $\dots$
- $z_{n,1}, \dots, z_{n,T}$
- $\dots$
- $z_{N,1}, \dots, z_{N,T}$

where

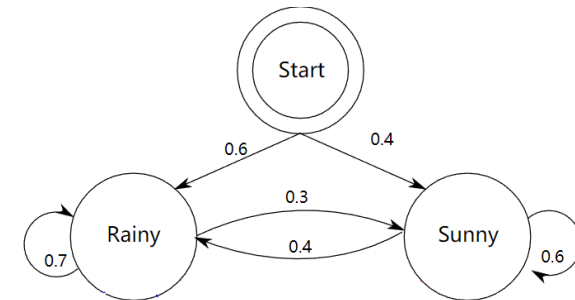
- for simplicity we assume each sequence has the same length  $T$
- lower case  $z_{n,t}$  represents the value of the random variable  $Z_{n,t}$

From these observations how do we *learn the model parameters*  $\theta := (\pi, A)$ ?

## Graph Representation

picture from Wikipedia

It is intuitive to represent a Markov model as a **graph**



## Finding the MLE

Same story, find the **MLE**. The log-likelihood of a sequence  $z_1, \dots, z_T$  is

$$\begin{aligned}
 & \ln p(Z_{1:T} = z_{1:T}; \theta) \\
 &= \sum_{t=1}^T \ln p(Z_t = z_t \mid Z_{1:t-1} = z_{1:t-1}; \theta) && (\text{always true}) \\
 &= \sum_{t=1}^T \ln p(Z_t = z_t \mid Z_{t-1} = z_{t-1}; \theta) && (\text{Markov property}) \\
 &= \ln \pi_{z_1} + \sum_{t=2}^T \ln a_{z_{t-1}, z_t} \\
 &= \sum_s \mathbb{I}[z_1 = s] \ln \pi_s + \sum_{s, s'} \left( \sum_{t=2}^T \mathbb{I}[z_{t-1} = s, z_t = s'] \right) \ln a_{s, s'}
 \end{aligned}$$

## Finding the MLE

So the MLE is

$$\begin{aligned} \operatorname{argmax}_{\theta} \ln p(\mathcal{D}; \theta) &= \operatorname{argmax}_{\theta} \sum_{n=1}^N \ln p(Z_{n,1:T} = z_{n,1:T}; \theta) \\ &= \operatorname{argmax}_{\pi, \mathbf{A}} \sum_s (\text{\#initial states with value } s) \ln \pi_s \\ &\quad + \sum_{s,s'} (\text{\#transitions from } s \text{ to } s') \ln a_{s,s'} \end{aligned}$$

subject to

$$\sum_{s'} a_{s,s'} = 1 \quad \text{and} \quad a_{s,s'} \geq 0 \quad \forall s \in [S]$$

We have seen this many times. The solution is:

$$\begin{aligned} \pi_s &\propto \text{\#initial states with value } s \\ a_{s,s'} &\propto \text{\#transitions from } s \text{ to } s' \end{aligned}$$

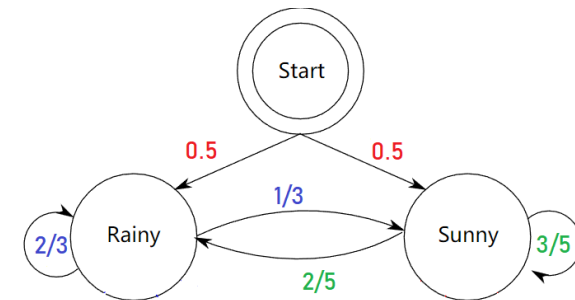
See also: [MLaPP 17.2.2] and <http://cs229.stanford.edu/section/cs229-hmm.pdf>

## Example

Suppose we observed the following 2 sequences of length 5

- sunny, sunny, rainy, rainy, rainy
- rainy, sunny, sunny, sunny, rainy

MLE is the following model



## Markov Model with outcomes

Now suppose each state  $Z_t$  also “emits” some **outcome**  $X_t \in [O]$  based on the following model

$$P(X_t = o \mid Z_t = s) = b_{s,o} \quad (\text{emission probability})$$

independent of anything else.

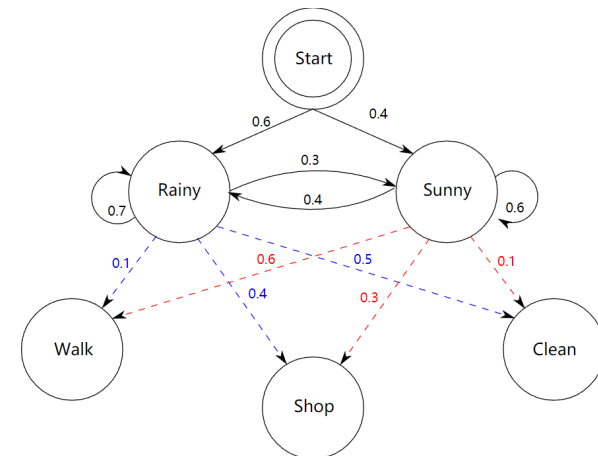
For example, in the language model,  $X_t$  is the speech signal for the underlying word  $Z_t$  (very useful for **speech recognition**).

Now the model parameters are  $(\{\pi_s\}, \{a_{s,s'}\}, \{b_{s,o}\}) = (\pi, \mathbf{A}, \mathbf{B})$ .

## Another example

picture from Wikipedia

On each day, we also observe **Bob's activity: walk, shop, or clean**, which only depends on the weather of that day.



## Joint likelihood

The joint log-likelihood of a **state-outcome sequence**  $z_1, x_1, \dots, z_T, x_T$  is

$$\begin{aligned}
 & \ln P(Z_{1:T} = z_{1:T}, X_{1:T} = x_{1:T}) \\
 &= \ln P(Z_{1:T} = z_{1:T}) + \ln P(X_{1:T} = x_{1:T} \mid Z_{1:T} = z_{1:T}) \quad (\text{always true}) \\
 &= \sum_{t=1}^T \ln P(Z_t = z_t \mid Z_{t-1} = z_{t-1}) + \sum_{t=1}^T \ln P(X_t = x_t \mid Z_t = z_t) \\
 & \hspace{15em} (\text{due to all the independence}) \\
 &= \ln \pi_{z_1} + \sum_{t=2}^T \ln a_{z_{t-1}, z_t} + \sum_{t=1}^T \ln b_{z_t, x_t}
 \end{aligned}$$

## Learning the model

However, *most often we do not observe the states!* Think about the speech recognition example.

This is called **Hidden Markov Model (HMM)**, widely used in practice

How to learn HMMs? **Roadmap:**

- first discuss how to **infer** when the model is known (key: **dynamic programming**)
- then discuss how to **learn** the model (key: **EM**)

## Learning the model

If we observe  $N$  state-outcome sequences:  $z_{n,1}, x_{n,1}, \dots, z_{n,T}, x_{n,T}$  for  $n = 1, \dots, N$ , the MLE can again be obtained in a similar way (verify yourself):

$$\begin{aligned}
 \pi_s &\propto \text{\#initial states with value } s \\
 a_{s,s'} &\propto \text{\#transitions from } s \text{ to } s' \\
 b_{s,o} &\propto \text{\#state-outcome pairs } (s, o)
 \end{aligned}$$