# §11 first page Markov chains

§12.1
Learning a
random process

## P(A and B and C) = P(A) P(B|A) P(C(B,A) Random process = P(c) P(A|c) P(B|A,c)a sequence $X_0, X_1, X_2, \dots$ of random variables, typically not independent $\Pr(x_0, x_1, \dots, x_n) = \Pr_{X_0}(x_0) \Pr_{X_1}(x_1 | x_0) \Pr_{X_2}(x_2 | x_0, x_1) \times \dots \times \Pr_{X_n}(x_n | x_0 \dots x_{n-1})$ by the chain rule for probability If we have a dataset of sequences, and we have a probability model (e.g. a RNN or a Transformer neural network) that computes $\Pr_{X_i}(x_i | x_0 \cdots x_{i-1})$ , then we can fit it using maximum likelihood estimation. Because $X_2$ is generaled based only on $X_1$ , $Pr_{X_2}(x_2|x_0, x_1) = Pr_{X_2}(x_2|x_1)$ , Markov chain a random process in which each $X_i$ is generated based **only** on the preceding value $X_{i-1}$ $X_0 \to X_1 \to X_2 \to \cdots$ $\Pr(x_0, x_1, \dots, x_n) = \Pr_{X_0}(x_0) \Pr_{X_1}(x_1 | x_0 (\Pr_{X_2}(x_2 | x_1))) \cdots \times \Pr_{X_n}(x_n | x_{n-1})$

# Applications of Markov chains: dynamical systems



Let  $X_t$  be the full state of the system at time t. We'd like to use historical data to learn the dynamics  $(X_t|X_{t-1} = x_{t-1})$ , so that we can simulate it.

# Applications of Markov chains: stable diffusion

Given an image, create a sequence with progressively more and more noise, until we get pure noise. Do this for many images, to create a training dataset of sequences.



Reverse the sequences. Train a Markov chain to learn the dynamics  $(X_t | X_{t-1} = x)$ .



If we apply these dynamics to a new pure-noise image, we will generate a novel image.



"Guess who's back" "Ta na na"

Eminem



#### Joseph Fourier (1768-1830)

## THÉORIE

ANALYTIQUE

# DE LA CHALEUR,

PAR M. FOURIER.



#### A PARIS,

CHEZ FIRMIN DIDOT, PÈRE ET FILS,

LIBRAIRES POUR LES MATHÉMATIQUES, L'ARCHITECTURE HYDRAULIQUE ET LA MARINE, RUE JACOB, Nº 24.

1822.

Two main results:

any function of a variable, whether continuous or discontinuous, can be expanded in a series of sines of multiples of the variable  $\rightarrow$  Fourier transform

partial differential equation for conductive diffusion of heat

 $\rightarrow$  heat diffusion depends on the distance from the heat source

Example 12.1.1: fitting a Markov model Let  $[x_0, x_1, ..., x_n]$  be a time series which we believe is generated by  $X_{i+1} = a + b X_i + N(0, \sigma^2).$ 

Estimate a, b, and  $\sigma$  using maximum likelihood estimation.



$$P_{r}(x_{o}x_{1}\cdots x_{n}) = P_{r}(x_{o}) \prod_{i=1}^{n} P_{r}(x_{i} | x_{i-1})$$

$$= P_{r}(x_{o}) \prod_{i=1}^{n} \frac{1}{12\pi\sigma^{2}} e^{-(x_{i}^{*} - (a+bx_{i-1}))^{2}/2\sigma^{2}}$$
since  $X_{i} \sim N(a+bX_{i}^{*}, \sigma^{2})$ ,
The question hells us
$$= ??? \prod_{i=1}^{n} \cdots \cdots$$
the dist-of  $X_{0}$ 

$$Pr(x_{o}, x_{i}, \dots, x_{n}) = \frac{777}{11} \frac{1}{\sqrt{2\pi\sigma^{2}}} e^{-(\pi i - (a + b\pi i))^{2}/2\sigma^{2}}$$

To fit our model, we need to maximize this expression over  $a, b, \sigma$ .

predictor	response
<i>x</i> <sub>0</sub>	<i>x</i> <sub>1</sub>
<i>x</i> <sub>1</sub>	<i>x</i> <sub>2</sub>
<i>x</i> <sub>2</sub>	<i>x</i> <sub>3</sub>
:	:
$x_{n-1}$	$x_n$

But this is exactly the same maximization as for the supervised learning task of predicting  $x_i$  given  $x_{i-1}$  using the model  $X_i \sim a + bx_{i-1} + N(0, \sigma^2)$ 

It's simple to fit using sklearn.

# Autoregressive modelling

This is a regression (i.e. supervised learning with numerical response). It's called 'auto' because we're predicting x using x itself as a predictor.

# §11.2 Calculations with Markov chains

There are three ways to specify a Markov chain model.



#### CAUSAL DIAGRAM

Each  $X_i$  is generated based only on the preceding state  $X_{i-1}$ :

$$X_1 \to X_2 \to X_3 \to \cdots$$

#### TRANSITION PROBABILITY MATRIX

$$P = \frac{\text{drizzle}}{\text{grey}} \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

 $P_{ij} = \mathbb{P}\left( \begin{array}{c} \text{next state} \\ \text{is } j \end{array} \middle| \begin{array}{c} \text{in state} \\ i \end{array} \right)$ 

If the state space is  $\mathbb{R}$  we can't write out the full matrix so we instead specify  $\Pr_{X_t}(x_t|X_{t-1} = x_{t-1})$ 

This is particularly used to describe diffusion models

Wait, diffusion like thermodynamics?!?!?





Horizontal steam boiler, Augsburg machines, early 19<sup>th</sup> century

Joseph Fourier (1768-1830)

Geille Sealt

Example 11.2.1 (Multi-step transition probabilities) If it's grey today, what's the chance of rain two days from now?

-

2



$$P(X_{2} = T \mid X_{0} = g)$$

$$r = rain
g = grey
d = drizzle.
$$\sum_{x} P(X_{2} = r \mid X_{1} = x, X_{0} = g) P(X_{1} = x \mid X_{0} = g)$$

$$Law of Total Probability
with baggage & X_{0} = g]$$

$$\sum_{x} P(X_{2} = T \mid X_{1} = x) P(X_{1} = x \mid X_{0} = g)$$
since  $X_{2}$  is generalized based only an  $X_{1}$ ,  
the stablestrime  $D$  is irrelevant one we  
know the close at time  $1$ .  
a.k.a. Memorylessness:  

$$\sum_{x} P_{xr} P_{gx} = \sum_{x} P_{gx} P_{xr} = [P^{2}]_{gr}$$$$

# Laws of probability that can help when working with Markov chains §11.2

Law of Total Probability $\mathbb{P}(A = a)$ $= \sum_{b} \mathbb{P}(A = a \mid B = b) \mathbb{P}(B = b)$	Law of Total Probability with baggage { $C = c$ } $\mathbb{P}(A = a \mid C = c)$ $= \sum_{b} \mathbb{P}(A = a \mid B = b, C = c) \mathbb{P}(B = b \mid C = c)$
Bayes's rule	<b>Bayes's rule</b> with baggage $\{C = c\}$
$\mathbb{P}(A = a \mid B = b)$	$\mathbb{P}(A = a \mid B = b, C = c)$
$= \frac{\mathbb{P}(A=a) \mathbb{P}(B=b A=a)}{\mathbb{P}(B=b)}$	$= \frac{\mathbb{P}(A = a   C = c) \mathbb{P}(B = b   A = a, C = c)}{\mathbb{P}(B = b   C = c)}$
<b>Definition of independence</b> If <i>A</i> and <i>B</i> are independent then $\mathbb{P}(A = a \mid B = b) = \mathbb{P}(A = a)$	<b>Definition of conditional independence</b> If <i>A</i> and <i>B</i> are conditionally independent given $\{C = c\}$ then $\mathbb{P}(A = a \mid B = b, C = c) = \mathbb{P}(A = a \mid C = c)$

#### Calculating with Markov Chains

The chain is memoryless  $X_0 \rightarrow X_1 \rightarrow \cdots$ i.e. each item is generated based only on the previous item Whenever we're doing calculations with Markov chains, we have to wrangle our expression into a form where we can use memorylessness (plus the transition probability matrix).

Often, this will involve conditioning using the Law of Total Probability.

#### The memorylessness theorem:

conditional on the present, the future is independent of the past.

$$\mathbb{P}(X_{3} = x_{3} | X_{2} = x_{2}, X_{1} = x_{1}, X_{0} = x_{0}) = \mathbb{P}(X_{3} = x_{3} | X_{2} = x_{2})$$

$$\mathbb{P}(X_{3} = x_{3} | X_{2} = x_{2}, X_{0} = x_{0}) = \mathbb{P}(X_{3} = x_{3} | X_{2} = x_{2})$$

$$\mathbb{P}(X_n = x_n \mid X_{n-1} = x_{n-1}, \dots, X_0 = x_0) = \mathbb{P}(X_n = x_n \mid X_{n-1} = x_{n-1})$$

# Technicalities (\*non-examinable)

Formally, a Markov chain is defined by specifying the form of its likelihood function:  $\forall x_0, \dots, x_n$  $\Pr(x_0, x_1, \dots, x_n) = \Pr_{X_0}(x_0) \Pr_{X_1}(x_1|x_0) \Pr_{X_2}(x_2|x_1) \times \dots \times \Pr_{X_n}(x_n|x_{n-1})$ 

From this, one can prove memorylessness results such as  $\Pr_{X_3}(x_3 \mid X_2 = x_2, X_1 = x_1, X_0 = x_0) = \Pr_{X_3}(x_3 \mid X_2 = x_2)$ 

and indeed the full memorylessness theorem.

If you're ever stuck trying to prove a result about Markov chains, and if you can't see a way to use memorylessness, try going back to basics in the form of the likelihood function.

#### Exercise

Given that yesterday was rain, and tomorrow is rain, what's the chance that today is drizzle?



today yesterday tomorrow  

$$P(X_{1}=d | X_{0}=r_{1} | X_{2}=r)$$
Bayes's rule:  $A \rightarrow B$ ,  
figure at A given B.  
What we're doing:  $X_{0} \rightarrow x_{1} \rightarrow x_{2}$   

$$P(X_{1}=x_{1} | X_{0}=x_{0}, X_{2}=x_{2})$$
What we're doing:  $X_{0} \rightarrow x_{1} \rightarrow x_{2}$   

$$P(X_{1}=x_{1} | X_{0}=x_{0}, X_{2}=x_{2})$$

$$P(X_{1}=x_{1}, X_{0}=x_{0}, X_{2}=x_{2})$$
by definition of analitianal probability  

$$P(X_{0}=x_{0}, X_{2}=x_{2})$$
nuwarabor =  $P(X_{0}=x_{0}, X_{1}=x_{1}, X_{2}=x_{2})$  by simple rewriting  
=  $P(X_{0}=x_{0}) P(X_{1}=x_{1} | X_{0}=x_{0}) P(X_{2}=x_{2} | X_{1}=x_{1})$   
using the generat form of libelihood for a Murthau chain  
(proved using the Chain Rule + Memory kesness)  
denomination =  $\sum_{y} P(X_{0}=x_{0}, X_{2}=x_{2}, X_{1}=y)$  by the Som Rule  
(a version of the Law of Tot Prob)  
=  $\sum_{y} P(X_{0}=x_{0}) P(X_{1}=y | X_{0}=x_{0}) P(X_{2}=x_{2} | X_{1}=y)$  above,  
=  $P(X_{0}=x_{0}) P_{x_{0}y}P_{y_{x_{2}}}$ 

$$= \frac{P_{x_{0}}x_{1}P_{x_{1}}x_{2}}{\sum_{y} P(X_{0}=x_{0}) P_{x_{0}}y_{y_{x_{2}}}}$$

# Why I'm excited about this sort of result (\* non-examinable)

In science, we don't just want to learn associations, we want to learn causal mechanisms.

For example, smoking is associated with getting cancer ... but perhaps smoking is protective against cancer, and the association is because of some hidden causal factor (e.g. genetics) that encourages smoking and also predisposes towards cancer.

In machine learning, we're often presented with a supervised learning task ("learn to predict y given  $x_1$  and  $x_2$ "), and we don't even think about the underlying mechanisms.

- If the causal mechanism is  $X_1 \rightarrow Y \rightarrow X_2$ , we can still train a supervised learning model to predict Y (as per the previous exercise)
- Open research question: how can we train ML systems to learn the causal mechanisms, rather than just associations?







# Hidden Markov models



For a hidden Markov model, the likelihood function  $Pr_{\underline{X}}(\underline{x})$  is nasty, and it's pretty much impossible to learn the model from  $\underline{x}$  data.

So why are hidden Markov models useful?

- Uber collects precise logs (both <u>z</u> and <u>x</u>) from a few drivers, so it can learn the full probability model for how <u>Z</u> and <u>X</u> are generated using straightforward supervised learning
- Then, for regular trips (only  $\underline{x}$  data available), they can infer the posterior  $(\underline{Z} | \underline{X} = \underline{x})$  using Bayes's rule
- (Alternatively, they can simply find the most likely Z<sub>T</sub> using the Viterbi algorithm)



