

Conditional Probability:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$\therefore P(B|A) = \frac{P(B \cap A)}{P(A)}$$

$$\Rightarrow P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Bayes' Theorem – False Positives

- You test positive in a test with a 5% false positive rate.
- What is the chance you have the disease?
- Suppose 1 in 100 have the illness. Does this matter?
- The test also has a false negative rate of 10%

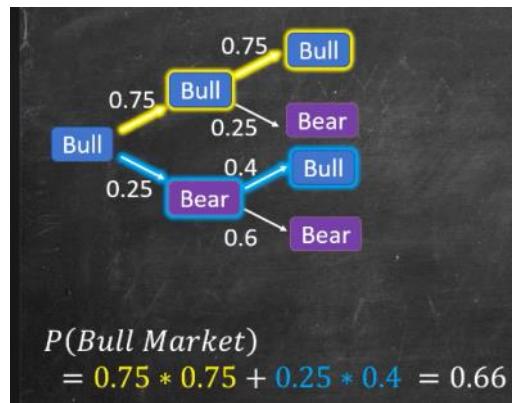
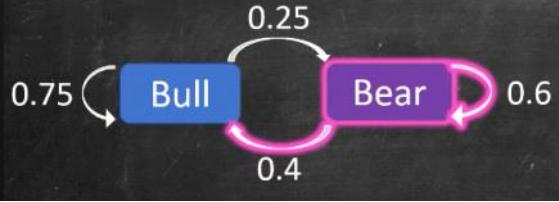


- The most important strength of bayes theorem is that , if you don't know one conditional probability , there is still the chance you know the other conditional probability.
- The rate at which something occur in a populace or a category, i.e rare or common affect the probability .
- As new knowledge comes, the probability would be updated .

A **Markov Chain** is a sequence of events where the probabilities of the future **ONLY** depend on the present

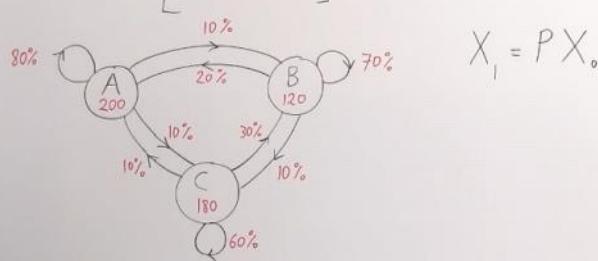
Example: Stock Market

- 75% a bull market followed by a bull market
- 60% a bear market followed by a bear market



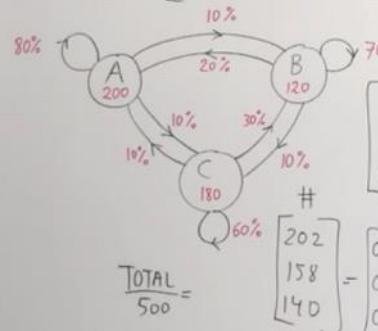
PROBABILITY AND STATISTICS : MARKOV CHAINS ① INTRODUCTION

$$[\text{NEXT STATE}] = \left[\begin{array}{c} \text{MATRIX OF} \\ \text{TRANSITION} \\ \text{PROBABILITIES} \end{array} \right] [\text{CURRENT STATE}]$$



PROBABILITY AND STATISTICS : MARKOV CHAINS (1) INTRODUCTION

$$[\text{NEXT STATE}] = [\text{MATRIX OF TRANSITION PROBABILITIES}] [\text{CURRENT STATE}]$$



$$X_1 = P X_0$$

$$\begin{aligned} \begin{bmatrix} A_1 \\ B_1 \\ C_1 \end{bmatrix} &= \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.7 & 0.3 \\ 0.1 & 0.1 & 0.6 \end{bmatrix} \begin{bmatrix} A \\ B \\ C \end{bmatrix} \\ &= \begin{bmatrix} 0.8(0.4) + 0.2(0.24) + 0.1(0.36) \\ 0.1(0.4) + 0.7(0.24) + 0.3(0.36) \\ 0.1(0.4) + 0.1(0.24) + 0.6(0.36) \end{bmatrix} \\ &= \begin{bmatrix} 0.404 \\ 0.316 \\ 0.280 \end{bmatrix} \end{aligned}$$



- The main idea behind these markov chain is that if **matrix of transition probabilities** is kept constant , after a large number of iteration the next state -> value matrix will start to converge.

PROBABILITY AND STATISTICS : MARKOV CHAINS (3) WHY THEY CALL IT CHAINS

$$\begin{aligned} X_1 &= P X_0 & X_2 &= P X_1 & X_3 &= P X_2 \dots \\ \begin{matrix} \text{NEXT STATE} \\ (\text{STATE } 1) \end{matrix} &\quad \begin{matrix} \text{INITIAL STATE} \\ \text{PROBABILITY MATRIX} \end{matrix} & & & & \\ \begin{bmatrix} 0.404 \\ 0.316 \\ 0.280 \end{bmatrix} &= \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.7 & 0.3 \\ 0.1 & 0.1 & 0.6 \end{bmatrix} \begin{bmatrix} 0.40 \\ 0.24 \\ 0.36 \end{bmatrix} & \begin{matrix} A = 200 \\ B = 120 \\ C = 180 \\ \hline 500 \end{matrix} & & & \\ X_2 &= \begin{bmatrix} 0.414 \\ 0.346 \\ 0.240 \end{bmatrix} &= \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.7 & 0.3 \\ 0.1 & 0.1 & 0.6 \end{bmatrix} \begin{bmatrix} 0.404 \\ 0.316 \\ 0.280 \end{bmatrix} &= \begin{bmatrix} 0.3(0.404) + 0.2(0.316) + 0.1(0.280) \\ 0.1(0.404) + 0.7(0.316) + 0.3(0.280) \\ 0.1(0.404) + 0.1(0.316) + 0.6(0.280) \end{bmatrix} & & \\ X_3 &= \begin{bmatrix} 0.417 \\ 0.343 \\ 0.240 \end{bmatrix} & & & & \end{aligned}$$

$$\begin{aligned} X_1 &= P X_0 \\ X_2 &= P X_1 \\ X_3 &= P X_2 \\ \text{OR} \\ X_2 &= P^2 X_0 \\ \text{AND} \\ X_3 &= P^3 X_0 \end{aligned}$$

- Once the convergence is reached , we have reached the stable distribution matrix and things stop to change if probability matrix is left unchanged.

PROBABILITY AND STATISTICS : MARKOV CHAINS (5) THE CHAIN CONTINUED

$$\begin{aligned} X_1 &= P X_0 = \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.7 & 0.3 \\ 0.1 & 0.1 & 0.6 \end{bmatrix} \begin{bmatrix} 0.404 \\ 0.316 \\ 0.280 \end{bmatrix} \\ X_2 &= P X_1 = \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.7 & 0.3 \\ 0.1 & 0.1 & 0.6 \end{bmatrix} \begin{bmatrix} 0.404 \\ 0.316 \\ 0.280 \end{bmatrix} \\ X_3 &= P X_2 = \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.7 & 0.3 \\ 0.1 & 0.1 & 0.6 \end{bmatrix} \begin{bmatrix} 0.414 \\ 0.346 \\ 0.240 \end{bmatrix} \\ X_4 &= P X_3 = \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.7 & 0.3 \\ 0.1 & 0.1 & 0.6 \end{bmatrix} \begin{bmatrix} 0.414 \\ 0.346 \\ 0.240 \end{bmatrix} \\ X_5 &= P X_4 = \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.7 & 0.3 \\ 0.1 & 0.1 & 0.6 \end{bmatrix} \begin{bmatrix} 0.424 \\ 0.356 \\ 0.220 \end{bmatrix} \end{aligned}$$

SAME

$$X_6 = P X_5 = \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.7 & 0.3 \\ 0.1 & 0.1 & 0.6 \end{bmatrix} \begin{bmatrix} 0.424 \\ 0.356 \\ 0.220 \end{bmatrix} = \begin{bmatrix} 0.432 \\ 0.358 \\ 0.210 \end{bmatrix}$$

$$X_7 = P X_6 = \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.7 & 0.3 \\ 0.1 & 0.1 & 0.6 \end{bmatrix} \begin{bmatrix} 0.432 \\ 0.358 \\ 0.210 \end{bmatrix} = \begin{bmatrix} 0.438 \\ 0.357 \\ 0.205 \end{bmatrix}$$

STABLE DISTRIBUTION MATRIX = \bar{X}

46% 45% 45%

A B C

200 160 140

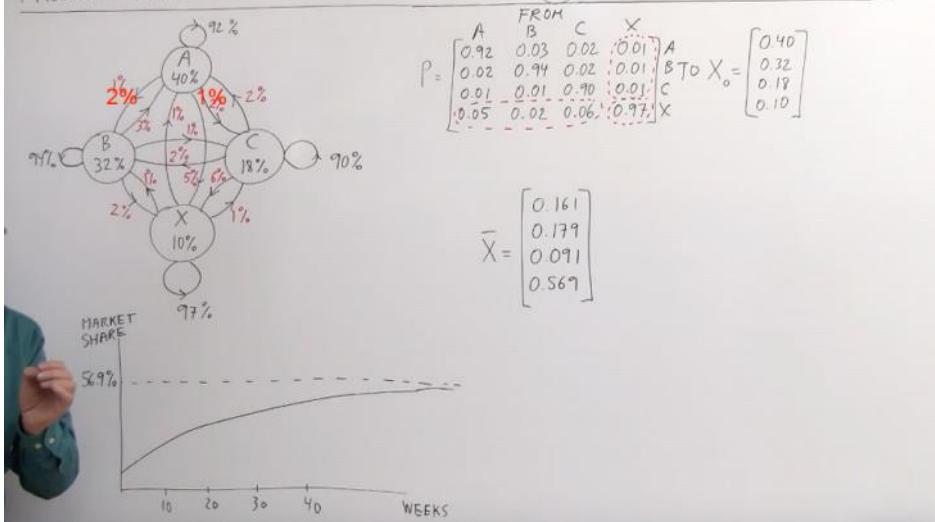
35% 35% 35%

B C

177 101

20% 20%

PROBABILITY AND STATISTICS : MARKOV CHAINS (6) MARKET PENETRATION



- The final state is kinda independent of initial state in this case , and this showcases the power of markov chains. If we have probability matrix , we can start from a given initial state and irrespective of it the final state will almost be same !!!

PROBABILITY AND STATISTICS : MARKOV CHAINS (7) POWERS OF THE PROBABILITY MATRIX

$$\begin{aligned}
 P &= \begin{bmatrix} A & B \\ 0.8 & 0.1 \\ 0.2 & 0.9 \end{bmatrix} A \rightarrow \\
 P^2 &= \begin{bmatrix} 0.8 & 0.1 \\ 0.2 & 0.9 \end{bmatrix} \begin{bmatrix} 0.8 & 0.1 \\ 0.2 & 0.9 \end{bmatrix} = \begin{bmatrix} 0.66 & 0.17 \\ 0.34 & 0.83 \end{bmatrix} \\
 P^3 &= \begin{bmatrix} 0.8 & 0.1 \\ 0.2 & 0.9 \end{bmatrix} \begin{bmatrix} 0.66 & 0.17 \\ 0.34 & 0.83 \end{bmatrix} = \begin{bmatrix} 0.562 & 0.219 \\ 0.438 & 0.781 \end{bmatrix} \\
 P^4 &= \begin{bmatrix} 0.8 & 0.1 \\ 0.2 & 0.9 \end{bmatrix} \begin{bmatrix} 0.562 & 0.219 \\ 0.438 & 0.781 \end{bmatrix} = \begin{bmatrix} 0.493 & 0.253 \\ 0.507 & 0.747 \end{bmatrix} \\
 P^5 &= \begin{bmatrix} 0.8 & 0.1 \\ 0.2 & 0.9 \end{bmatrix} \begin{bmatrix} 0.493 & 0.253 \\ 0.507 & 0.747 \end{bmatrix} = \begin{bmatrix} 0.445 & 0.277 \\ 0.555 & 0.723 \end{bmatrix} \\
 P^6 &= \begin{bmatrix} 0.8 & 0.1 \\ 0.2 & 0.9 \end{bmatrix} \begin{bmatrix} 0.445 & 0.277 \\ 0.555 & 0.723 \end{bmatrix} = \begin{bmatrix} 0.412 & 0.294 \\ 0.538 & 0.706 \end{bmatrix} \\
 P^\infty &= \bar{P} = \begin{bmatrix} 0.333 & 0.333 \\ 0.667 & 0.667 \end{bmatrix} = \begin{bmatrix} \frac{1}{3} & \frac{1}{3} \\ \frac{2}{3} & \frac{2}{3} \end{bmatrix} \\
 \bar{P} X_0 &= ? = \bar{X} = \text{STABLE DISTRIBUTION MATRIX} \\
 \begin{bmatrix} \frac{1}{3} & \frac{1}{3} \\ \frac{2}{3} & \frac{2}{3} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} &= \begin{bmatrix} \frac{1}{3} \\ \frac{2}{3} \end{bmatrix} \\
 \begin{bmatrix} \frac{1}{3} & \frac{1}{3} \\ \frac{2}{3} & \frac{2}{3} \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} &= \begin{bmatrix} \frac{1}{3} \\ \frac{2}{3} \end{bmatrix} \\
 \begin{bmatrix} \frac{1}{3} & \frac{1}{3} \\ \frac{2}{3} & \frac{2}{3} \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} &= \begin{bmatrix} \frac{1}{3} \\ \frac{2}{3} \end{bmatrix}
 \end{aligned}$$

- Markov chain only works when the probability matrix is stochastic and regular in nature.

PROBABILITY AND STATISTICS : MARKOV CHAINS (8) STOCHASTIC MATRIX

$$\begin{aligned}
 P &= \begin{bmatrix} A & B \\ 0.25 & 0.50 \\ 0.75 & 0.50 \end{bmatrix} A \rightarrow X_0 = \begin{bmatrix} A \\ B \end{bmatrix} \quad X_1 = X_0 P \\
 &\downarrow \text{ADD TO 1} \quad \therefore \text{STOCHASTIC}
 \end{aligned}$$

$$\begin{aligned}
 P &= \begin{bmatrix} A & B \\ 0.25 & 0.75 \\ 0.50 & 0.50 \end{bmatrix} \quad X_0 = \begin{bmatrix} A & B \end{bmatrix} \\
 &\xrightarrow{\text{ADD TO 1}} \quad \therefore \text{STOCHASTIC}
 \end{aligned}$$

$$\begin{aligned}
 &\text{NOT STOCHASTIC} \\
 &\text{FROM} \\
 &\downarrow \begin{bmatrix} 0.9 & 0.1 \\ 0.2 & 0.8 \end{bmatrix} \text{TO} \\
 &\neq 1
 \end{aligned}$$

$$\begin{aligned}
 &\text{NOT STOCHASTIC} \\
 &\text{TO} \\
 &\text{From} \begin{bmatrix} 0.9 & 0.2 \\ 0.3 & 0.9 \end{bmatrix} \\
 &\xrightarrow{\quad} \neq 1
 \end{aligned}$$

PROBABILITY AND STATISTICS : MARKOV CHAINS (9) REGULAR MATRIX

STOCHASTIC $P = \begin{bmatrix} 0.25 & 0.50 \\ 0.75 & 0.50 \end{bmatrix}$ FROM TO

A STOCHASTIC MATRIX IS ALSO A REGULAR MATRIX
IF P^n ($n > 1$) HAS ONLY NON-ZERO POSITIVE ENTRIES

$$P = \begin{bmatrix} 0.25 & 0.50 \\ 0.75 & 0.50 \end{bmatrix}$$

$$P^2 = \begin{bmatrix} 0.25 & 0.50 \\ 0.75 & 0.50 \end{bmatrix} \begin{bmatrix} 0.25 & 0.50 \\ 0.75 & 0.50 \end{bmatrix}$$

$$= \begin{bmatrix} 0.438 & 0.375 \\ 0.562 & 0.625 \end{bmatrix}$$

∴ REGULAR MATRIX

$$P = \begin{bmatrix} 1 & 0.5 \\ 0 & 0.5 \end{bmatrix}$$

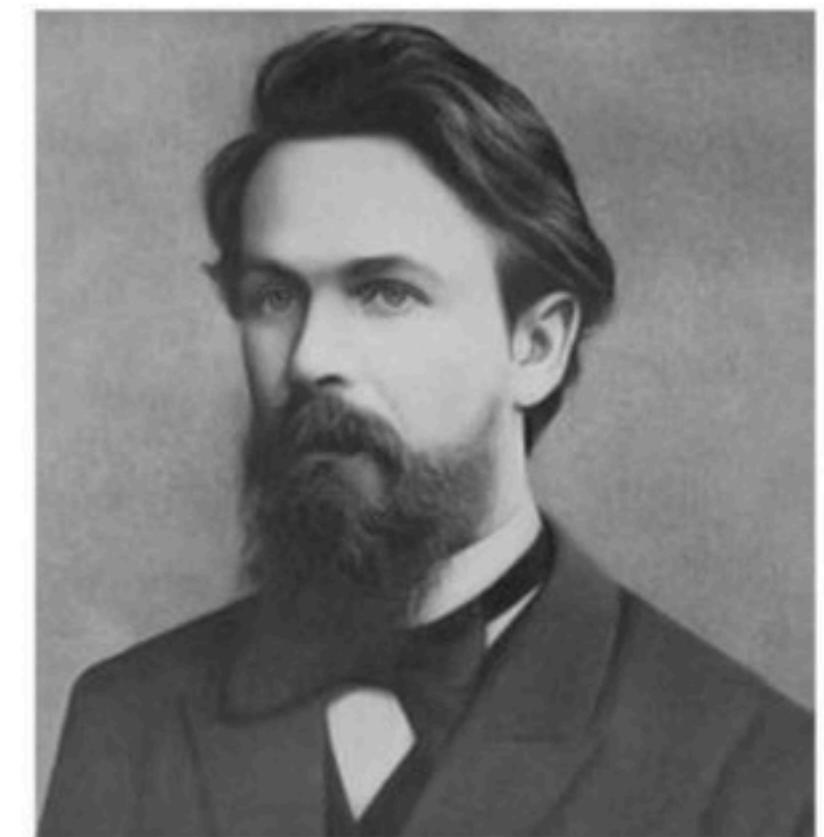
$$P^2 = \begin{bmatrix} 1 & 0.5 \\ 0 & 0.5 \end{bmatrix} \begin{bmatrix} 1 & 0.5 \\ 0 & 0.5 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0.75 \\ 0 & 0.25 \end{bmatrix}$$



Markov Models - Section Introduction

- Markov models are everywhere!
- Finance: the basis for the famous Black-Scholes formula
- Reinforcement Learning: Markov Decision Process (MDP)
- Hidden Markov Model (speech recognition, computational biology)
- Markov Chain Monte Carlo (MCMC): numerical approximation

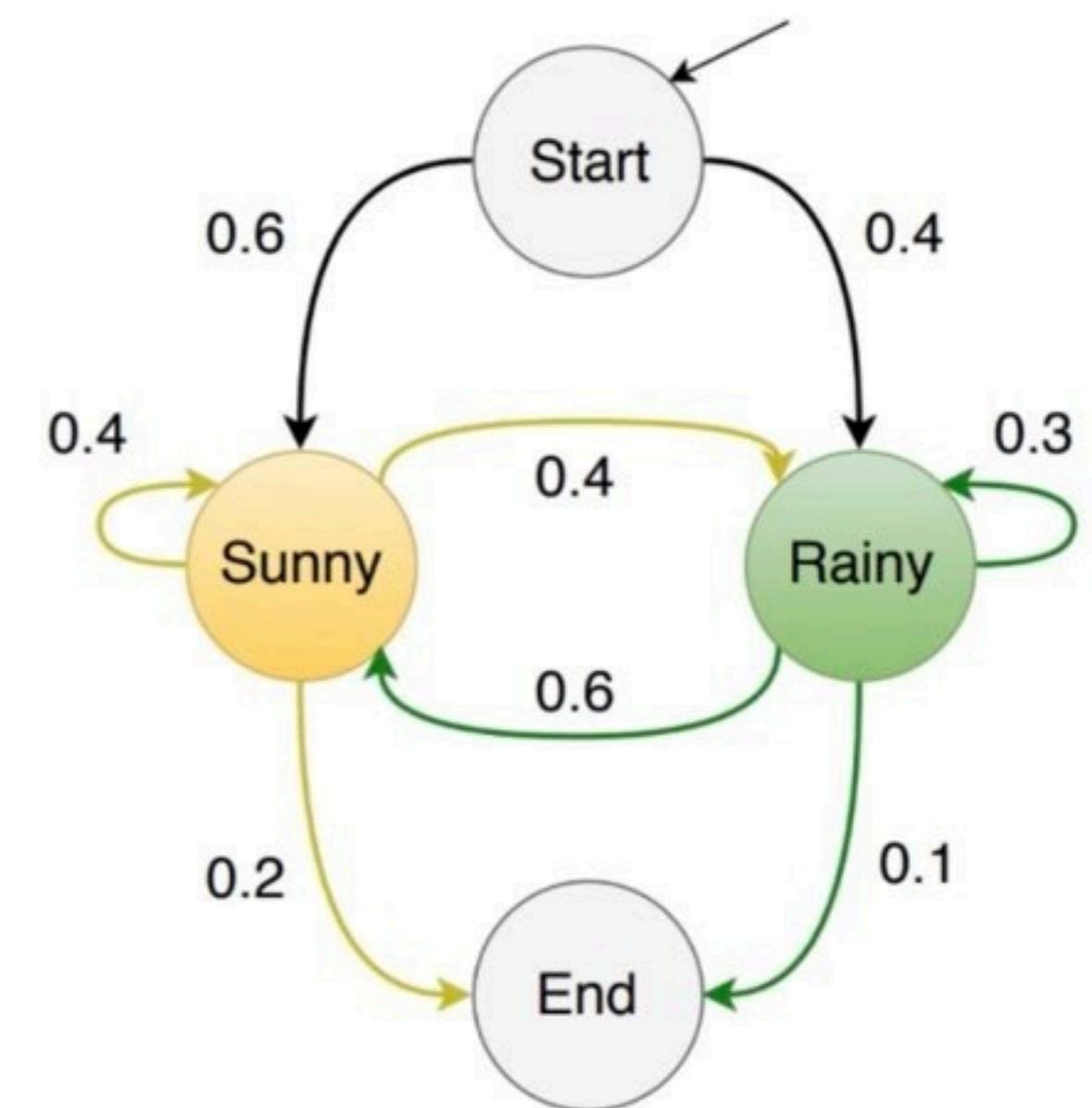


Odemys

Another major application of the Markov model is Markov Chain Monte Carlo, also known as EMC, which

Section Outline

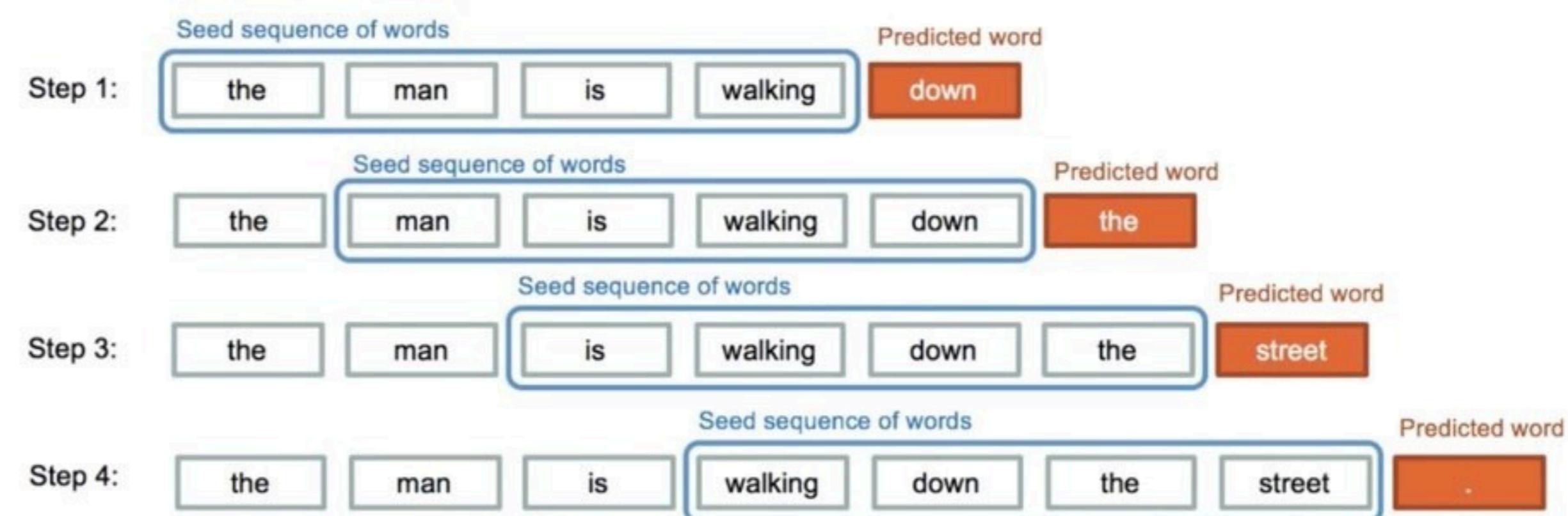
- Markov property (requires knowledge of probability)
- Markov model, state transition matrix
- "Training" a Markov model (i.e. "learning")



We'll also look at how to train a Markov model.

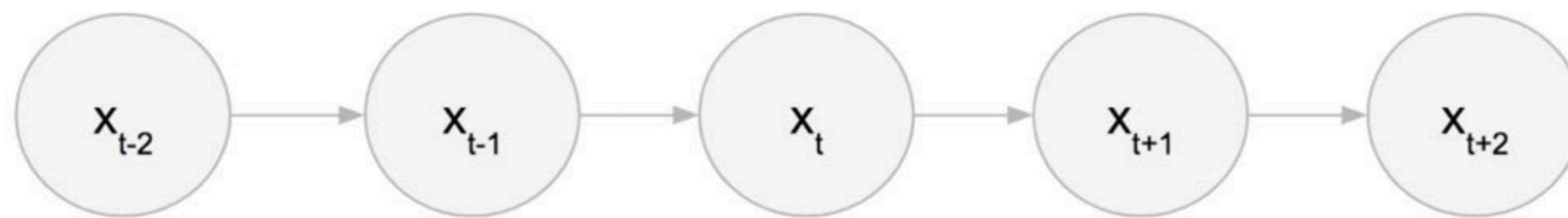
Applications

- Building a text classifier using Bayes rule + Generating poetry
- These are 2 fundamentally different ways to apply ML
- Supervised: predict a target from a dataset of inputs + labels
- Unsupervised: learn structure of data, create new examples having the same structure



The Markov Property

- The Markov property is a very restrictive assumption on the dependency structure of the joint distribution (arrows denote dependency)



- This implies that the joint distribution takes the form:

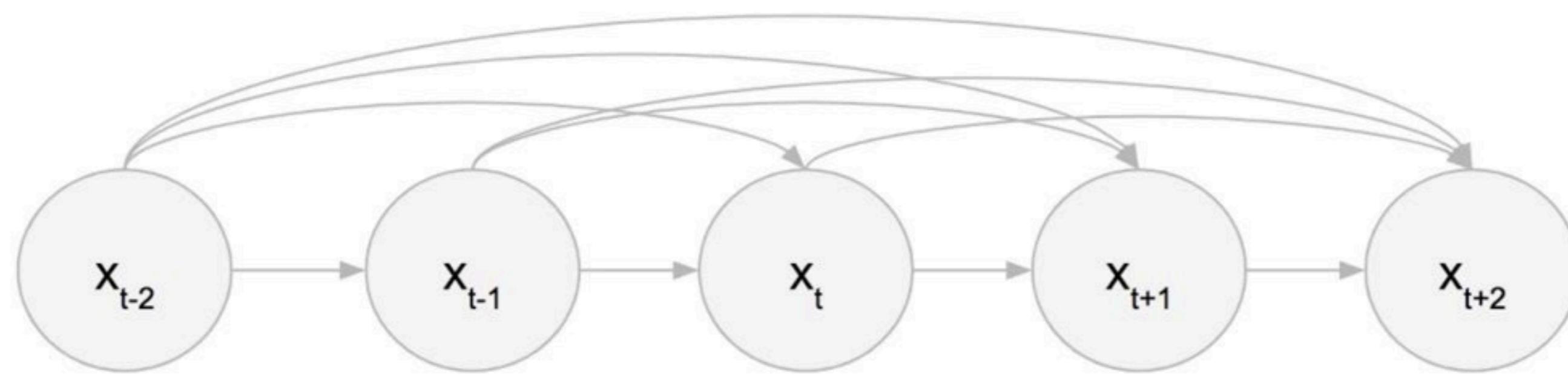
$$p(x_1, \dots, x_T) = p(x_1) \prod_{t=2}^T p(x_t | x_{t-1})$$

an x one times p of x three given x two and so on, up to T given x t minus y.

Chain Rule of Probability

- T-variable expansion (general case)
- In general, each term depends on **all** preceding terms

$$p(x_1, \dots, x_T) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2)\dots p(x_T | x_{T-1}, \dots, x_1)$$



So you see, we've kind of learned about the mark of property backwards.

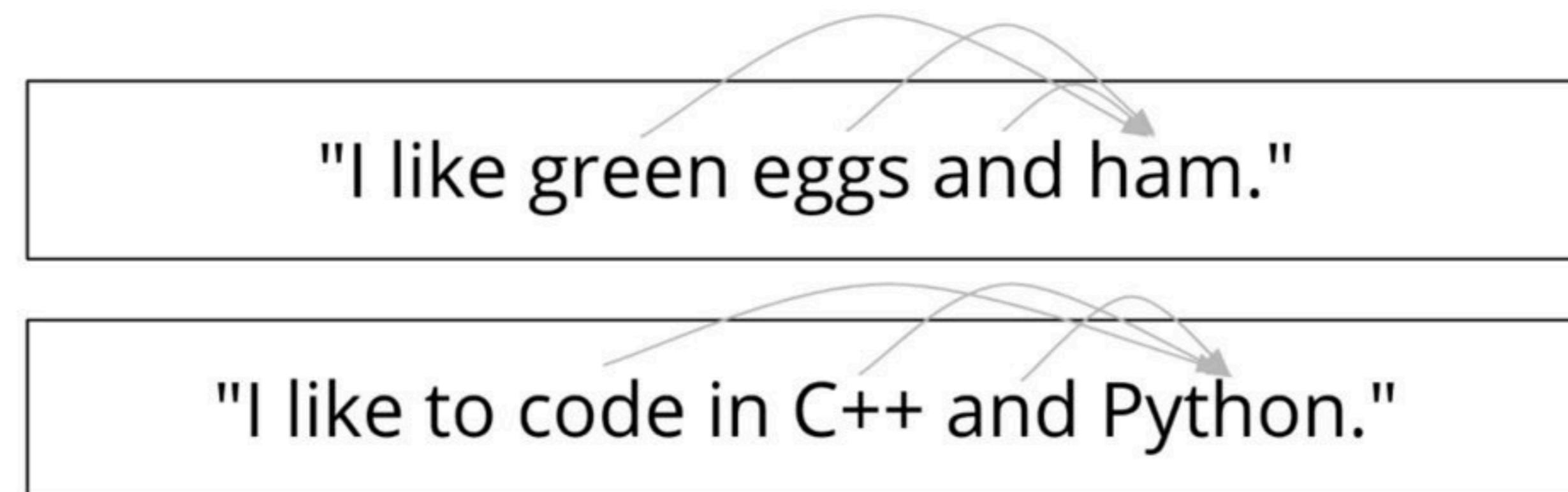
Why is the Markov Property Useful?

- Consider a Markov model of the English language
- Suppose we choose 2000 of the most common words (though some estimate that each individual uses 20k - 30k words)
- Suppose we want to model sentences of length 10 (not unreasonable)
- Our model is $p(x_{10} \mid x_9, x_8, \dots, x_1)$
- Each "x" has 2000 possible values (words)
- How many probabilities to estimate? $2000 \times 2000 \times 2000 \dots = 2000^{10}$

In order to estimate all these numbers, we must have a sufficient amount of data to learn from as another

The Markov "Assumption"

- We assume the Markov property holds, even when it does not

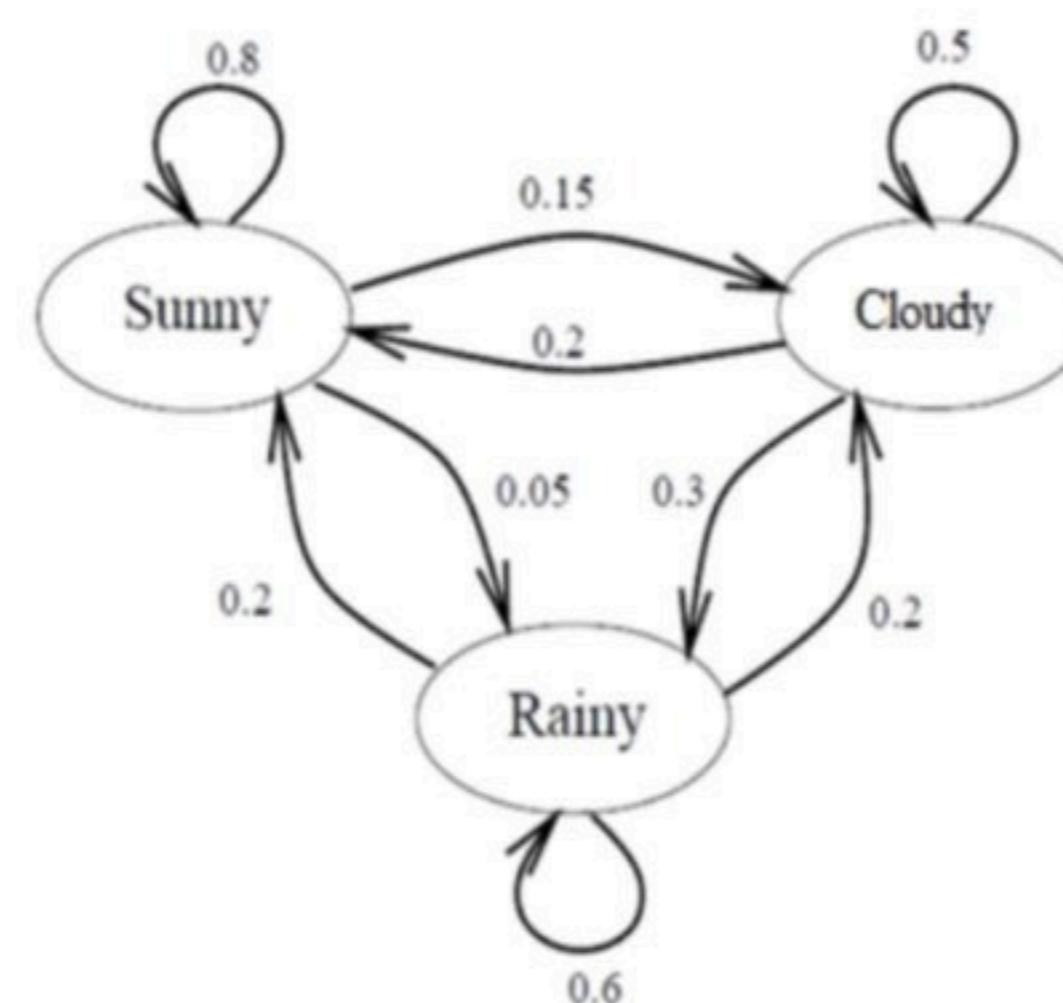


- Do "ham" and "Python" **only** depend on "and"? No!

We can also see that the word python becomes more likely because code in C++ appeared earlier in the

The Markov Model

- In this section, we'll consider sequences of **categorical symbols**, which we'll refer to as "states"
- Example weather: {sunny, rainy, cloudy}
- Example parts-of-speech tags: {noun, verb, adjective, ...}



verb, adjective, and so forth.

Notation

- We'll use s for "state" (other sources might use z or x)

$s(t) = s_t = \text{state at time } t$

- Time will also be discrete ($t = 1, 2, \dots$) (e.g. no such thing as $t = 1.5$)
- We will **number** the states from $1, 2, \dots, M$
- M = total number of possible states (e.g. $M = 3$ for sunny/rainy/cloudy)
- We use i or j to **index** the state space

$p(s_t = i)$ means : probability that state at time t is ' i '

So for example, press of T equals I translates to the probability that the state at time t is equal

State Distribution

- We have $p(s_t = 1), p(s_t = 2), \dots, p(s_t = M)$
- This is M probability values - together they form a **distribution**
- Ex. "What is the probability that it will be rainy on Sunday?"
Ans: $p(s_{\text{sunday}} = \text{rainy})$

$p(s_t)$ = state distribution (length M vector)

State Transitions

$$p(s_t = j \mid s_{t-1} = i)$$

- "Probability that state at time t is j, given that the state at time t-1 was i"
- How many of these probability values exist? (for all i and j)
- Since both i and j can take any value from 1...M, there are M^2 values

Well, we must have a value for every possible i.e from one up to M and every possible j from one up

State Transition Matrix

- Wait a minute... what happened to t?

$$A_{ij} = p(s_t = j \mid s_{t-1} = i), \forall i = 1 \dots M, j = 1 \dots M$$

- In general, we could have $A_{ij}(t)$
- When A doesn't depend on t: **time-homogeneous Markov process**

...

Probability of a Sequence

$$p(s_{1...T}) = p(s_1)p(s_2|s_1)p(s_3|s_1, s_2)\dots$$

$$p(s_{1...T}) = p(s_1) \prod_{t=2}^T p(s_t|s_{t-1})$$

Markov Property

$$p(s_{1...T}) = \pi_{s_1} \prod_{t=2}^T A_{s_{t-1}, s_t}$$

Start thinking about how you'd implement this in Python code

Thus you should consider how you would write some code to carry out this computation.

Training a Markov Model

- We're not ready to rigorously derive this (yet), but the intuition is simple
- Suppose we flip a coin a bunch of times - how do we estimate $p(\text{heads})$?

This answers the question: "How frequent is the word 'cat' in the English language?"

- That's just the binary case. What if we have $M > 2$?

$$p("cat") \approx \frac{\text{count}("cat")}{\text{total word count}}$$

...

Estimating A and π ("Training")

Note: the 'hat'
means 'estimate'

N = number of
sequences in dataset

$$\hat{\pi}_i = \frac{\text{count}(s_1 = i)}{N}$$

$$\hat{A}_{ij} = \frac{\text{count}(i \rightarrow j)}{\text{count}(i)}$$

Example:
 $\text{count("the cat") / count("the")}$

We would like to know what is the probability of seeing the word cat following the word the.

Markov Model Summary

- What we want to do: model a sequence of states
- State transition matrix (A), initial state distribution (π)
- Task #1: find the probability of a sequence
- Task #2: given a dataset, find A and π
- Task #2 is called "learning" or "training"
- Start considering how to implement these in computer code

Now, again, it's always useful to think about how you would implement this in computer code.

Small Modifications

- In this lecture, we will:
 - 1) Modify how we estimate A and π
 - 2) Compute the probability of a sequence

Previously, we used
"maximum likelihood
estimation"

Probability of a Sequence

$$p(s_1 \dots s_T) = \pi_{s_1} \prod_{t=2}^T A_{s_{t-1}, s_t}$$

- Only involves multiplication
- What if one of the values is 0? Transition never appeared in train set
- The whole thing becomes 0!

This will make the entire expression zero, although this seems like it might be the right answer.

Add-One Smoothing

- Give a small probability to *every* possible transition
- Add a "fake count" of 1 for each (i, j) transition

$$\hat{A}_{ij} = \frac{\text{count}(i \rightarrow j) + 1}{\text{count}(i) + M}$$

Adding M to the denominator ensures that each row of A sums to 1

Since we had a one to each of the impossible values, we must also add em to the denominator as well.

Add-One Smoothing

- We can do this for the initial state distribution too

$$\hat{\pi}_i = \frac{count(s_1=i) + 1}{N + M}$$

Add-Epsilon Smoothing

- We can make it more smooth ($\varepsilon > 1$) or less smooth ($\varepsilon < 1$)

$$\hat{A}_{ij} = \frac{\text{count}(i \rightarrow j) + \varepsilon}{\text{count}(i) + \varepsilon M}$$

$$\hat{\pi}_i = \frac{\text{count}(s_1 = i) + \varepsilon}{N + \varepsilon M}$$

denominator again by the same logic we use before we can conclude that this results in each row of AIG,

Computing the Probability of a Sequence

$$p(s_1 \dots T) = \pi_{s_1} \prod_{t=2}^T A_{s_{t-1}, s_t}$$

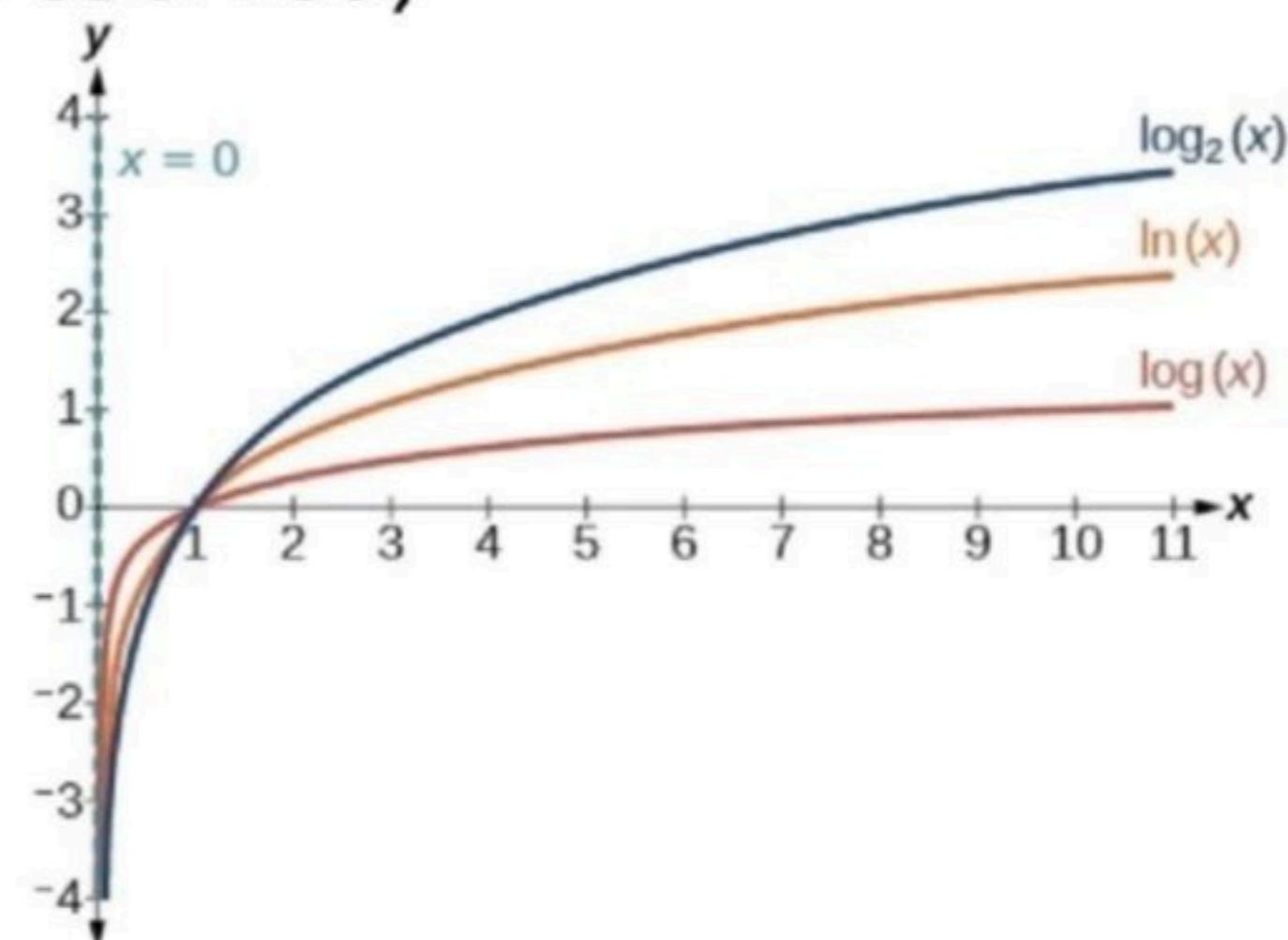
- This involves **multiplying** many small numbers together
- Common to use 20k - 50k vocabulary size in English
- As you multiply small numbers together, they approach 0! ($0.1 * 0.1 = 0.01$)
- Computers don't have infinite precision; eventually it just rounds to 0
- Problematic: what if we want to compare 2 sequences?

Try it yourself in wait until you encounter it to convince yourself that it can and will happen.

Working with Log Probabilities

- Solution: compute log probabilities instead
- We don't need the actual probability **value**, since what we usually want to do is **compare** (e.g. is one sequence more likely than another?)
- It works because the log is a monotonically increasing function
- If $A > B$, then $\log(A) > \log(B)$ (i.e. the "winner" is preserved)
- Does it work? $\log_{10}(10^{-10}) = -10$
- Does it work? $\log_{10}(10^{-100}) = -100$

Now you are approaching the limit
of what the computer can represent!



Working with Log Probabilities

$$\log p(s_{1:T}) = \log \pi_{s_1} + \sum_{t=2}^T \log A_{s_{t-1}, s_t}$$

- Since $\log(AB) = \log(A) + \log(B)$
- If we don't need the probabilities, we can store only the logged values
- Note: do not compute the product and **then** take the log
- Also note: addition (+) is more efficient (faster) than multiplication (*)

computing a product, as you recall, adding is much faster than doing multiplication.

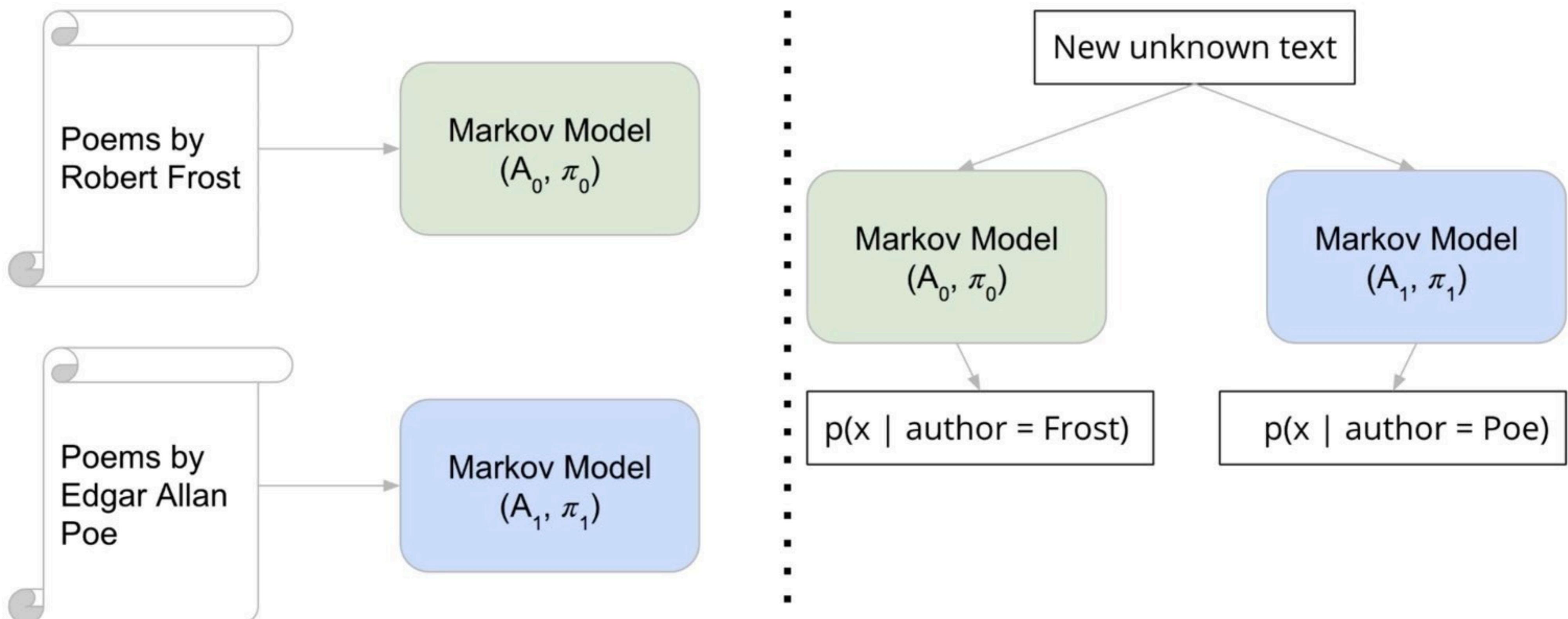
Supervised or unsupervised?

- Text classification is an example of supervised learning, but Markov models are unsupervised (training data is just text; no labels)
- Answer: We must apply Bayes' rule (and build a **Bayes classifier**)

$$p(y | x) = \frac{p(x | y)p(y)}{p(x)}$$

In fact, what we are really building in this lecture is a Bayes classifier.

Bayes Classifier



Applying Bayes' Rule

$$p(\text{author} \mid \text{poem}) = \frac{p(\text{poem} \mid \text{author}) p(\text{author})}{p(\text{poem})}$$

So this is the basic form of what we need, but there are several ways we can simplify this expression

Simplifying the decision rule

$$k^* = \arg \max_k \frac{p(poem | author = k) p(author = k)}{p(poem)}$$

$$k^* = \arg \max_k p(poem | author = k) p(author = k)$$

$$k^* = \arg \max_k \log p(poem | author = k) + \log p(author = k)$$

So what we end up with is that we would like the ARG mass of log of palm giving author plus
log p of

Maximum Likelihood

- We can simply take the argmax of the likelihood if the prior is uniform (i.e. we have no reason to believe, given no other info, that one author is more likely than another)

$$k^* = \arg \max_k \log p(poem | author = k)$$

If $p(\text{author})$ is uniform

Recap

- We train a separate Markov model for each class
- Each model gives us $p(x | \text{class} = k)$ for all k
- General form of decision rule using Bayes' rule: $k^* = \operatorname{argmax}_k p(\text{class} = k | x)$
- Posterior can be simplified since we don't need its actual value
- Maximum a posteriori (MAP): $k^* = \operatorname{argmax}_k \log(p(x | \text{class} = k)) + \log(p(\text{class} = k))$
- Maximum likelihood: $k^* = \operatorname{argmax}_k \log(p(x | \text{class} = k))$

Text Classifier Exercise Prompt

- You'll be given poems by 2 authors: Edgar Allan Poe and Robert Frost
- URLs will be in the coming notebooks
 - Copy the URLs from there, don't try to search / type by hand (yes, really)
 - But don't look at the rest of the notebook, of course!
- Build a classifier that can distinguish between the 2 authors
- Compute train and test accuracy
- Check for class imbalance, compute F1-score if imbalanced



Details

- Some details are specific to my implementation
 - If you had other ideas, they may not apply to you (use your own judgment)
-
- Loop through each file, save each line to a list (one line == one sample)
 - Save the labels too
 - Train-test split
 - Create a mapping from unique word to unique integer index
 - Loop through data, tokenize each line (string split should suffice)
 - Assign each unique word a unique integer index ("mapping" == "dict")
 - Create a special index for unknown words (words in the test set but not the train set)

This is because there may be words in the test set that do not appear in the train set.

Details

- Convert each line of text (the samples) into integer lists
- Train a Markov model for each class (Edgar Allan Poe / Robert Frost)
- Don't forget to use smoothing (e.g. add-one smoothing)
- Consider whether you need A and π , or $\log(A)$ and $\log(\pi)$
- For Bayes' rule, compute the priors: $p(\text{class} = k)$
- Now you have everything you need to make a prediction!



Odemy

At this point, you should have everything you need to make a prediction.

Details

- Write a function to compute the posterior for each class, given an input
- Take the argmax over the posteriors to get the predicted class
- Make predictions for both train and test sets
- Compute accuracy for train/test
- Check for class imbalance
- If imbalanced, check confusion matrix and F1-score



Furthermore, you may want to compute a metric like the F1 score, which takes into account any imbalance