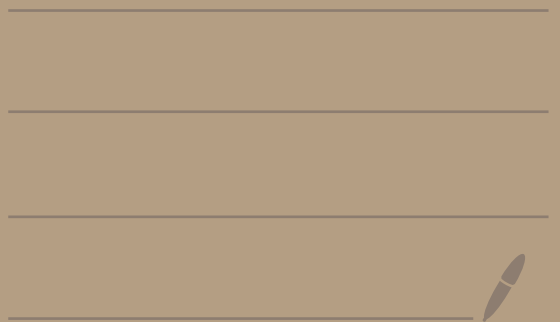


D3 Markov Models



Markov Property

Sequence: $\{x_1, x_2, x_3, \dots\}$ timeseries / words

Probability of sequence $p(x_1, x_2, x_3, \dots)$

Predict: $p(x_T | x_{T-1}, x_{T-2}, \dots)$

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

$p(x_t | x_{t-1}, x_{t-2}, \dots) = p(x_t | x_{t-1})$ nothing before only previous element influences
 \Rightarrow independent of x_{t-2}, x_{t-3}, \dots

What if MP not true?

$$p(x_1, x_2) = p(x_1) p(x_2 | x_1)$$

$$p(x_1, x_2, x_3) = p(x_1) \underbrace{p(x_2 | x_1) p(x_3 | x_2, x_1)}_{p(x_2, x_3 | x_1)}$$

Chain rule of probability

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

Example

Find 10th word, given the first 9

$p(x_{10} | x_1 \dots x_9)$ x has 2000 possible words

Estimate 2000^{10} possibilities

Markov Assumption

eggs and ham C++ and Python

word depends only on "and"

The Markov Model

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

$t = 1, 2, 3, 4 \dots M$

$P(s_t = i)$ probab. state at $t = i$

State distribution

$P(s_t = 1) P(s_t = 2) \dots P(s_t = M)$

State transitions

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

prob. state at t is j , given that state at $t-1$ was i

State transition matrix

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

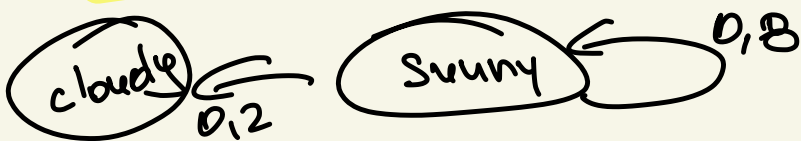
first index = prev state

second index = next state

in general $A_{ij}(t)$ possible, but assume

time homogeneous Markov process: A_{ij} static

State transition Diagram



Initial State

initial state distribution

$\pi_i = P(s_1 = i)$ for $i = 1 \dots M$

π is $M \times 1$ $A = M \times M$

Probability of sequence

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

state transition

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

Training Markov Model

$$p(\text{heads}) \approx \frac{\text{count}(\text{H})}{\text{total tosses}}$$

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

Estimating A/π

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

seq. started with i
number of sequences

estimate with hist

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

Probability smoothing / log prob.

A/π = max likelihood estimate

$$P(s_1 \dots \tau) = \pi_{s_1} \prod_{t=2}^{\tau} A_{s_t s_{t-1}}$$

\downarrow
= 0 since not appeared in training set (only in testing)

Add one smoothing

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

\rightarrow each row of $A_{ij} = 1$

$$\frac{1}{N} = \frac{\text{count}(s_1=i) + 1}{N + M}$$

ϵ can be $\epsilon \geq 1$ and $M = \epsilon M$

more smooth $\epsilon > 1$

less $\epsilon < 1$

Compute prob sequence

A/π are very small (20-50k engl. words)
multiply gets closer to zero
compare 2 improbable sentences

Working with log probabilities

$$A > B \quad \log(A) > \log(B)$$

$$\log_{10}(10^{-100}) = -100$$

$$\log P(s_1 \dots \tau) = \log \pi_{s_1} + \sum_{t=2}^{\tau} \log A_{s_t s_{t-1}} \quad \rightarrow \text{adding faster}$$

Text Classifier

Poem \rightarrow classify \rightarrow Poe
 \rightarrow Frost

email \rightarrow spam

movie review \rightarrow sentiment

Supervised or unsupervised

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

Frost $\rightarrow A_0 \pi_0 \rightarrow p(x|\text{author}=\text{Frost})$

Poe $\rightarrow A_1 \pi_1 \rightarrow p(x|\text{author}=\text{Poe})$

Apply Bayes' Rule

$p(\text{poem}|\text{author})$ but want $p(\text{author}|\text{poem})$

$$\Rightarrow k^* = \underset{k}{\operatorname{argmax}} p(\text{class}=k|x)$$

$x = \text{poem}$
 $k = \text{author}$

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

$$k^* = \underset{k}{\operatorname{argmax}} p(\text{poem}|\text{author}=k) p(\text{author}=k)$$

$\operatorname{argmax}(\log + \log)$

$p(\text{author})$ is uniform \Rightarrow const so

Max Likely hood $k^* = \underset{k}{\operatorname{argmax}} p(\text{poem}|\text{author}=k)$

Markov Models to generate Text

Classifying text \rightarrow supervised

Generating \rightarrow unsupervised (no labels)

Discriminative

$$p(y|x)$$

no need to apply
Bayes rule

Generative

$$p(x|y)$$

we find $p(y|x)$
later using Bayes rule

Sampling

$$N(0,1)$$

\downarrow \searrow
mean variance

Problems of Markov assumption

next word only depends on previous word

I made myself a peanut butter sandwich.

I will go and see her myself

Second order Markov Model:

$$p(s_t | s_{t-1}, s_{t-2} \dots) = p(s_t | s_{t-1}, s_{t-2})$$

2nd Order

$$A_{ijk} = p(s_t = k | s_{t-1} = j, s_{t-2} = i)$$

3D array (M^3 shape)

exp growth with number of past states

Full Model

$\pi_i = p(s_1 = i)$ for first word

$A^{(1)}_{ij}$ for second word

$A^{(2)}_{ijk}$ for every word after that

Article spinning

search engines detect duplicates

select a word to replace via a suggestion

N-Gram Approach

First order Markov Model for words

$$p(w_t | w_{t-1}) = \frac{\text{count}(w_{t-1} \rightarrow w_t)}{\text{count}(w_{t-1})}$$

$$p(w_t | w_{t-1}, w_{t-2}) = \frac{\text{count}(w_{t-2} \rightarrow w_{t-1} \rightarrow w_t)}{\text{count}(w_{t-2} \rightarrow w_{t-1})}$$

Predicting the middle word

$$p(w_t | w_{t-1}, w_{t+1}) = \frac{\text{count}(w_{t-1} \rightarrow w_t \rightarrow w_{t+1})}{\text{count}(w_{t-1} \rightarrow \text{ANY} \rightarrow w_{t+1})}$$

max likelihood est.

Production \rightarrow Began \rightarrow To

- \ Capacity /
- \ closer /
- \ continued /
- \ Facilities /

noun/verb/adj
may not belong