

Lesson 1: Word prediction via Markov chains [I]

(or writing like your favorite author)

In this lesson of introductory nature, we start by something easy: word prediction. The idea is to see a text as a sequence of random words. These are not independent & the simplest tool to model weak dependence is Markov chains. In fact, Andrei Andreevich Markov (1856-1922) designed Markov chains with the example of TEXT in mind - only that states were just "vowel" & "consonant".

[I] Formal definition:

Let $E = \{1, 2, \dots, n\}$ be a finite set, "the state space"

$P = (P_{ij})_{1 \leq i, j \leq n}$ is a stochastic matrix

$$\text{ie } \left(\sum_j P_{ij} = 1 \quad \forall i \iff P \mathbf{1} = \mathbf{1} \right)$$

Then the Markov chain with $\left\{ \begin{array}{l} \text{state space } E \\ \text{transition matrix } P \end{array} \right.$ is the E -valued stochastic process $(X_n)_{n \in \mathbb{N}}$ defined by

$$\bullet X_0 \in E$$

$$\bullet P_{x,y} = \mathbb{P}(X_{n+1} = y \mid X_n = x) \\ = \mathbb{P}(X_{n+1} = y \mid X_n = x, X_{n-1}, \dots)$$

This is a generalization of independent sequences II

- here dependence is only through the previous state.

The following important theorem generalizes the law of large numbers.

Thm [Birkhoff's ergodic Thm for Markov chains]

If X is irreducible

$$\text{Then } \frac{1}{T} \sum_{t=0}^T f(X_t) \xrightarrow[T \rightarrow \infty]{\text{a.s.}} \mathbb{E}_{x \sim \pi} (f(x)) = \sum_{x \in E} f(x) \pi(x)$$

where π is a probability measure

called the invariant measure, solution to $\pi P = \pi$

II The historical example Markov Andrei was in disagreement with Pavel Nekrasov that independence is necessary for the (weak) law of large numbers.

To that endeavor, in a paper written in 1913, Markov

chose a sequence of $T = 20000$ ~~words~~ letters

from Pushkin's "Eugene Onegin". He obtained

the Markov chain with state space $E = \{V, C\}$

$$P = \begin{matrix} & \begin{matrix} V & C \end{matrix} \\ \begin{matrix} V \\ C \end{matrix} & \begin{pmatrix} 0,128 & 0,872 \\ 0,663 & 0,337 \end{pmatrix} \end{matrix}$$

Computing π from $\pi P = \pi$, one finds $\pi = (0,432, 0,568)$

This generalizes easily to the problem of word prediction by setting

$E = \{ \text{words appearing in a text} \}$

$P =$ Transition matrix for the Markov chain whose realization is our text.

↳ How to estimate P ? How did Markov do?

Take a long text (say $T = 20000$ letters like Markov)

& compute $\hat{P}_T = \begin{matrix} v \\ c \end{matrix} \begin{pmatrix} \hat{\alpha}_T & 1 - \hat{\alpha}_T \\ 1 - \hat{\beta}_T & \hat{\beta}_T \end{pmatrix}$

where $\begin{cases} \hat{\alpha}_T = \frac{\# \{ 1 \leq t \leq T-1 \mid X_t = X_{t+1} = v \}}{\# \{ 1 \leq t \leq T-1 \mid X_t = v \}} \\ \hat{\beta}_T = \frac{\# \{ \dots \mid X_t = X_{t+1} = c \}}{\# \{ \dots \mid X_t = c \}} \end{cases}$

↳ why does it work?

Only the denominators are amenable to Birkhoff's theorem

as $\frac{1}{T} \# \{ 1 \leq t \leq T-1 \mid X_t = x \} = \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{1}_{\{X_t = x\}} \rightarrow \pi(x)$

Numerators can be treated with an extended version of Birkhoff. However, we shall prefer a very generic method in statistics & Machine learning

III - Maximum Likelihood Estimation.

IV

Proposition: $[\hat{P}_T]_{ij} := \frac{m_{ij}}{\sum_{j'} m_{ij'}}$ & $m_{ij} := \sum_{1 \leq t \leq T-1} \mathbb{1}_{(X_t, X_{t+1}) = (i, j)}$
 defines a matrix \hat{P}_T ; which is the consistent
 MLE estimator for P

Proof: • Deriving Likelihood & CV.

Suppose we observe the realization $(x_0, x_1, x_2, \dots, x_T)$,
 and x_0 is deterministically chosen initial state.
 Then the likelihood is the probability of our realization:

$$\begin{aligned} L &= \mathbb{P}(X_0 = x_0, X_1 = x_1, \dots, X_m = x_m) \\ &= \mathbb{P}(X_1 = x_1 | X_0 = x_0) \dots \mathbb{P}(X_m = x_m | X_{m-1} = x_{m-1}) \mathbb{P}(X_0 = x_0) \\ &= \mathbb{P}(X_0 = x_0) \prod_{t=0}^{T-1} P_{x_t, x_{t+1}} = \mathbb{P}(X_0 = x_0) \prod_{i, j=1}^m P_{ij}^{m_{ij}} \end{aligned}$$

\Rightarrow

MLE:

$$\begin{aligned} \hat{P}_T &:= \underset{P \text{ stochastic matrix}}{\text{Argmax}} \log L(P) \\ &= \underset{P}{\text{Argmax}} \sum_{i, j=1}^m m_{ij} \log P_{ij} \end{aligned}$$

Thanks to general theorems of existence of
 MLE estimators (for example C. Paganini's course)
 $\hat{P}_T \xrightarrow{T \rightarrow \infty} P$ stochastic matrix.

• Computing \hat{P}_T :

IV

⚠ The set of stochastic matrices has n constraints

→ Method of Lagrange Multipliers

$$\mathcal{L}(P, \lambda) := \sum_{i,j=1}^m m_{ij} \log P_{ij} - \sum_{i=1}^m \lambda_i \left(\sum_{j=1}^m P_{ij} - 1 \right)$$

for $\begin{cases} P \text{ matrix} \\ \lambda \in \mathbb{R}^n \end{cases}$

(P, λ) is critical

$$\Leftrightarrow \begin{cases} \forall i, \frac{\partial \mathcal{L}}{\partial \lambda_i} = 0 \\ \forall i, j, \frac{\partial \mathcal{L}}{\partial P_{ij}} = 0 \end{cases} \Leftrightarrow \begin{cases} \sum_{j=1}^m P_{ij} = 1 & \forall i \\ \forall i, j, 0 = \frac{m_{ij}}{P_{ij}} - \lambda_i \end{cases}$$

$$\Leftrightarrow \begin{cases} \forall i, j, P_{ij} = \frac{m_{ij}}{\lambda_i} \\ 1 = \sum_{j=1}^m P_{ij} = \left(\sum_{j=1}^m m_{ij} \right) / \lambda_i \end{cases}$$

Hence
the
result

□