



Stochastic Methods Rooted in Statistical Mechanics Markov Chains: Definitions of States & Transitions Stationarity & Ergodic Probabilities Global & Detailed Balance Equations

> Prof. Vasilis Maglaris <u>maglaris@netmode.ntua.gr</u> <u>www.netmode.ntua.gr</u> Room 002, New ECE Building Tuesday March 4, 2025

STOCHASTIC PROCESSES & OPTIMIZATION IN MACHINE LEARNING Statistical Mechanics & Machine Learning

- Long-term statistics of sample element x(n) are related to macroscopic concepts of mechanical physics reaching dynamic equilibrium under a given temperature. There is an analogy of terms such as entropy (disorder) with states and control parameters of Machine Learning systems
- Statistical Mechanics offers inference models based on measured (and assumed) statistics of input elements x(n) that enable self-organization algorithms (e.g. via Unsupervised Learning for data compression, classification into clusters, forming context-aware maps, correction of corrupted data) and generation of elements distributed according to probabilistic assumptions of the sample space (Statistical Sampling GenAI)
- The features of sample elements are encoded into random variables that constitute the *m* coordinates of input sample vectors selected from the environment (sample space).
 Analogies with statistical mechanics may lead to better understanding the impact of the *m* features of a large number of vectors used to train and validate a machine learning system
- Ground-braking application: **Boltzmann Machine** (*Hinton Sejnowski*, 1983) used for processing and generation of images inspired by statistical mechanics models (named after physicist and philosopher *Ludwig Boltzmann*, 1844-1906)

Statistical Mechanics: Gibbs Distribution, Partition Function, Entropy

Thermal Equilibrium of Physical System with many Degrees of Freedom

 A physical system with many degrees of freedom reaches dynamic equilibrium at absolute temperature T in state i of energy E_i with probability p_i (frequency of occurrence of i) given by the Gibbs Distribution (1902) or Boltzmann (1868):

$$p_i \propto \exp\left(-\frac{E_i}{T}\right), \quad \frac{p_i}{p_j} = \exp\left(-\frac{E_i - E_j}{T}\right), \qquad p_i \ge 0, \quad \sum_i p_i = 1$$

• States *i* of low energy *E_i* tend to occur more often in equilibrium with probabilities:

$$p_i = \frac{1}{Z} \exp\left(-\frac{E_i}{T}\right), \ Z = \sum_i \exp\left(-\frac{E_i}{T}\right)$$

Normalization Constant (Zustadsumme) Z: Partition

- The energy of the station is $E_i = -T \log(Zp_i)$ with an average $\langle E \rangle = \sum_i p_i E_i$
- The total Free Energy F (Helmholtz Free Energy) is:

$$F = -T \log Z \implies \langle E \rangle - F = -T \sum_{i} p_i \log p_i$$

• The entropy of the system is:

 $H \triangleq -\sum_{i} p_{i} \log p_{i} \implies \langle E \rangle -F = TH \notin F = \langle E \rangle -TH$ **Principle of Minimal Free Energy (Landau & Lifshitz**, 1980) In thermal equilibrium, H tends to its maximum, F to its minimum and states follow the **Gibbs Distribution**

Stochastic Processes - Time Series, Markov Property



- Stochastic Process of State X(t) with transitions from time τ to time t of the same sample element (time series outcome, trajectory) with probability $P\{[X(t) = a] | [X(\tau) = b]\}$
- Stochastic Process in Discrete Transition Times of State X_n ≜ X (n × Δt) with transitions from τ = (k × Δt) to t = (n × Δt) of the same sample element (discrete time series outcome, trajectory) with probability P(X_n = a|X_k = b)

Markov Property in Discrete Time Transition Processes

A discrete time Stochastic Process exhibits the Markov property if transition probabilities from state $X_n \rightarrow X_{n+1}$ do not depend on its past state evolution $\{X_1, X_2, ..., X_{n-1}\}$ $P(X_{n+1} = x_{n+1} | X_n = x_n, ..., X_1 = x_1\} = P(X_{n+1} = x_{n+1} | X_n = x_n)$

STOCHASTIC PROCESSES & OPTIMIZATION IN MACHINE LEARNING Markov Processes of Discrete State - Markov Chains

Consider a Markov Process of **Discrete States** $X_n = i$, referred to as **Markov Chain**, with **Discrete Transition Times** for $(X_n = i) \rightarrow (X_{n+1} = j)$ at instant *n*, independent of the past

Assume that transition probabilities in one step are constant and independent of instant *n*:

$$p_{ij} = P(X_{n+1} = j | X_n = i) \ge 0, \quad \sum_j p_{ij} = 1 \quad \forall i$$

If $i \le K$ the transition probabilities are elements of the transition matrix **P** with row elements adding to one: $\sum_{i} p_{ij} = 1$ (stochastic matrix):

$$\mathbf{P} = \begin{bmatrix} p_{11} & \cdots & p_{1K} \\ \vdots & \ddots & \vdots \\ p_{K1} & \cdots & p_{KK} \end{bmatrix}$$

The transition probabilities in *m* steps are $p_{ij}^{(m)} = P(X_{n+m} = j | X_n = i)$, m = 1,2,... are given by the **Chapman-Kolmogorov** identity:

$$p_{ij}^{(m+1)} = \sum_{k} p_{ik}^{(m)} p_{kj}$$
 kal $p_{ij}^{(m+n)} = \sum_{k} p_{ik}^{(m)} p_{kj}^{(n)}$, $m, n = 1, 2, ...$

STOCHASTIC PROCESSES & OPTIMIZATION IN MACHINE LEARNING Definitions of Markov Chain States

- **Recurrent States**: The process visits these states infinitely often in an infinite time-horizon
- **Transient States**: After a finite number of transition steps, the process stops visiting these states
- **Periodicity**: If all recurrent states are grouped in *d* disjoint subsets $S_1, S_2, ..., S_d$ with transitions allowed from a given subset to a next, then visits into subsets S_k occur periodically, every *d* transitions:

$$\text{If } i \in S_k, \, p_{ij} > 0 \ \Rightarrow \begin{cases} j \in S_{k+1} & \text{for } k = 1, \dots, d-1 \\ j \in S_1 & \text{for } k = d \end{cases}$$

- Irreducible Markov Chains: Two states communicate i ↔ j if the probability of reaching each other in finite number of steps is non-zero
 If i ↔ j & i ↔ k ⇒ i ↔ k (transitivity)
 If all states communicate the chain is Irreducible
- Classes: States can be classified in subsets. Open classes allow exits to a different class.
 Closed classes do not allow external transitions and (possible after a transient interval) the process is limited within this subset. If there is a single closed calss, the chain is Irreducible

d=3

6

Summary of Markov Chain State Classification



STOCHASTIC PROCESSES & OPTIMIZATION IN MACHINE LEARNING Transition Properties of Markov Chain (1/2)

• Mean Recurrence Time

In a Markov Chain (MC) the Mean Recurrence Time $E[T_i(k)]$ is defined as the average number of transitions (steps) for a Recurrent State *i* to cycle (return to itself), if it accomplished k-1 cycles from *i* to *i*. The relative occurrence of state *i* is proportional to the Steady State Probabilities:

$$\pi_i = \frac{1}{\mathrm{E}[T_i(k)]}$$

If $E[T_i(k)] < \infty$ then $\pi_i > 0$ and *i* is **Positive Recurrent**. Else $\pi_i = 0$ and *i* is **Null Recurrent**

• Steady State Probabilities, Ergodicity After *l* cycles into a *Positive Recurrent State i*, the proportion of steps (time) the MC resides in state *i* (*Sojourn Time*) is:

$$\nu_i(l) = \frac{l}{\sum_{k=1}^l T_i(k)}$$

The **Recurrence Times** $T_i(k)$ form a series of **Independent Identically Distributed** (IID) random variables. For $l \rightarrow \infty$ the proportions of transition steps the **MC** resides in state *i* approach to the **Steady State Probabilities** π_i :

$$\lim_{l\to\infty}\nu_i(l)=\pi_i,\ i=1,2,\ldots,K$$

The limit formula defines *Ergodicity* of state *i* (Ensemble Average = Time Average)

Transition Properties of Markov Chain (2/2)

MC State Transition Diagram: 3 State Example



STOCHASTIC PROCESSES & OPTIMIZATION IN MACHINE LEARNING Stationary Distribution of MC - Ergodic Probabilities

- Steady-state probabilities π_i of a **MC** are evaluated by measuring the relative frequency that state *i* occurs in an infinite (very large) repetition of sample outcomes
- Equivalently and for *irreducible* states steady-state probabilities can be deduced (measured) from the *trajectory* of a *single* outcome (time series) as the fraction of time that the process spends in *i* over an infinite (very lengthy) horizon.
- Such MC is defined as *ergodic* and the steady-state probabilities π_i as *ergodic probabilities*. An *irreducible* non-periodic MC is always ergodic

Convergence of State Probabilities to Ergodic Distributions

MC state probabilities of X_i , i = 1, 2, ..., K at transition step n = 0, 1, 2... define a $(1 \times K)$ vector $\pi^{(n)}$ evolving according to the $(K \times K)$ **stochastic transition matrix P** starting from an initial conditions $\pi^{(0)}$:

$$\boldsymbol{\pi}^{(n)} = \left[\pi_1^{(n)} \ \pi_2^{(n)} \dots \pi_K^{(n)} \right], \ \boldsymbol{\pi}^{(n)} = \boldsymbol{\pi}^{(n-1)} \mathbf{P} = \boldsymbol{\pi}^{(n-2)} \mathbf{P}^2 = \dots = \boldsymbol{\pi}^{(0)} \mathbf{P}^n$$

The ergodic probabilities $\pi = [\pi_1 \ \pi_2 \ ... \ \pi_K]$ are defined as the limit $\pi^{(n)}$ with $n \to \infty$:

$$\lim_{n \to \infty} \pi^{(n)} = \pi^{(0)} \times \lim_{n \to \infty} \mathbf{P}^n = \pi^{(0)} \begin{bmatrix} \pi_1 & \dots & \pi_K \\ \vdots & \ddots & \vdots \\ \pi_1 & \dots & \pi_K \end{bmatrix} = \pi^{(0)} \begin{bmatrix} \pi \\ \vdots \\ \pi \end{bmatrix} = \sum_{j=1}^K \pi_j^{(0)} \times \pi = 1 \times \pi = \pi$$

Hence $\pi = \lim_{n \to \infty} \pi^{(n)}$ is independent of the initial condition $\pi^{(0)}$ and can be evaluated by

solving the linear system of K-1 **Global Equilibrium Equations**, with an additional normalization linear equation to enforce linear independence of K equations;

 $\pi_j = \sum_{i=1}^{K} \pi_i p_{ij}$ j = 1, 2, ..., K or $\pi = \pi \mathbf{P}$ and $\sum_{j=1}^{K} \pi_j = 1$

STOCHASTIC PROCESSES & OPTIMIZATION IN MACHINE LEARNING Examples of Ergodic Markov Chains

State Transition Diagrams

States are denoted as circles x_1, x_2, \dots and transitions $x_i \rightarrow x_j$ as arrows with probabilities p_{ij}



Calculation of Ergodic Probabilities



Steady-state transition linear system

 $\pi_{1} = \frac{1}{3}\pi_{2} + \frac{3}{4}\pi_{3}$ $\pi_{2} = \frac{1}{6}\pi_{2} + \frac{1}{4}\pi_{3}$ $\pi_{3} = \pi_{1} + \frac{1}{2}\pi_{2}$ $\pi_{1} + \pi_{2} + \pi_{3} = 1 \text{ (normalization)}$ $\pi_{1} = 0.3953, \qquad \pi_{2} = 0.1395, \quad \pi_{3} = 0.4652$

STOCHASTIC PROCESSES & OPTIMIZATION IN MACHINE LEARNING Detailed Balanced Equations, Time Reversibility

- In thermal equilibrium a system that reaches state probabilities Gibbs $\pi_i = \frac{1}{Z} exp\left(-\frac{E_i}{T}\right)$ with Z the Partition Function, state transitions $i \rightarrow j$ are balanced in the long term by reverse transitions $j \rightarrow i$
- In this case *Detailed Balance Equations* hold with the ergodic probabilities π_i satisfying *Global Balance Equations* as expected for all Markov Chains (*MC*):



$$\sum_{i=1}^{K} \pi_{i} p_{ij} = \sum_{i=1}^{K} \left(\frac{\pi_{i}}{\pi_{j}} p_{ij} \right) \pi_{j} = \sum_{i=1}^{K} p_{ji} \pi_{i} = \pi_{i}$$

- The **Detailed Balance Equations** (unlike the **Global Balance Equations**) do not hold for all **MC**'s. An **MC** that abides by them is referred to as **Time Reversible** with forward transition probabilities p_{ij} and backwards with transition probabilities $\hat{p}_{ji} = \frac{\pi_i}{\pi_i} p_{ij}$
- The ergodic state probabilities π_i are much simpler to analyze if time-reversibility holds (e.g. for queuing networks under conditions yielding *product-form* state probabilities). Unfortunately, this is not true in many realistic models but holds for systems in thermal equilibrium that converge to *Gibbs* statistics (as in many *Machine Learning* systems)