



LINMA2470 Stochastic Modelling

SIMON DESMIDT

Academic year 2024-2025 - Q2

Contents

1	Reminders	2
1.1	General properties of probability	2
1.2	Expectation and variance	2
1.3	Law of large numbers	2
1.4	Central limit theorem	3
1.5	Exponential distribution	3
2	Poisson Processes	4
2.1	Distribution of $N(t)$	5
2.2	Non-homogenous Poisson processes	6
2.3	Bernoulli process approximation	6
2.4	Classification of queueing systems	7
3	Renewal Processes	8
3.1	Strong law of large numbers	8
3.2	Central limit theorem	8
3.3	Stopping time	9
3.4	Blackwell's renewal theorem	9
3.5	Little's Law	12
4	Finite State Markov Chains	14
4.1	Definitions	14
4.2	Transition probabilities	15
4.3	Markov chains with rewards	15
4.4	Markov decision processes	16
4.5	Dynamic programming algorithm for the stationary optimal policy . . .	17
5	Markov Decision Processes and Reinforcement Learning	18
5.1	MDP and Policies	18
5.2	Optimizing Policies	20
5.3	Value Iteration Algorithm	21

Reminders

1.1 General properties of probability

- $P[A \cup B] = P[A] + P[B] - P[A \cap B]$;
- $P[A|B] = \frac{P[A \cap B]}{P[B]} = \frac{P[AB]}{P[B]}$;
- A and B are independent iff $P[AB] = P[A]P[B] \implies P[A|B] = P[A]$;
- $P[X \leq x] = F_X(x)$ is the distribution function, i.e. a monotone increasing function of x going from 0 to 1 when x goes from $-\infty$ to $+\infty$.
- Its derivative is the density function $f_X(x)$ such that $f_X(x)\delta \approx P[x \leq X \leq x + \delta]$ for an infinitesimal δ .
- A random variable X is said to be memoryless if $\forall t, x > 0, P[X > t + x | X > t] = P[X > x]$.
- Markov inequality (for a nonnegative random variable): $P[Y \geq y] \leq \frac{\mathbb{E}[Y]}{y}$;
- Chebyshev inequality: $P[|Z - \mathbb{E}[Z]| \geq \varepsilon] \leq \frac{\sigma_Z^2}{\varepsilon^2}$;

1.2 Expectation and variance

- For a discrete random variable, $\mathbb{E}[X] = \sum_{n=-\infty}^{\infty} nP[X = n]$;
- For a continuous random variable, $\mathbb{E}[X] = \int_{-\infty}^{\infty} xf_X(x)dx$;
- $\mathbb{E}[X] = \int_0^{\infty} (1 - F_X(x))dx$.
- $Var[X] = \sigma_X^2 = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$;

1.3 Law of large numbers

Let X_1, \dots, X_n be a series of independent and uniformly distributed (IID) random variables with expectation \bar{X} and finite variance σ_X^2 . Let $S_n = X_1 + \dots + X_n$. Then,

- Weak version:

$$\lim_{n \rightarrow \infty} P \left[\left| \frac{S_n}{n} - \bar{X} \right| \geq \varepsilon \right] = 0 \quad (1.1)$$

- Strong version:

$$\lim_{n \rightarrow \infty} P \left[\sup_{m \geq n} \left(\frac{S_m}{m} - \bar{X} \right) > \varepsilon \right] = 0 \iff \lim_{n \rightarrow \infty} \frac{S_n}{n} = X \quad \text{with probability 1} \quad (1.2)$$

1.4 Central limit theorem

$$\lim_{n \rightarrow \infty} P \left[\frac{S_n - n\bar{X}}{\sqrt{n}\sigma} \leq y \right] = \int_{-\infty}^y \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \quad (1.3)$$

1.5 Exponential distribution

- $f_X(x) = \lambda e^{-\lambda x}$, for $x \geq 0$;
- $F_X(x) = 1 - e^{-\lambda x}$, for $x \geq 0$;
- $\mathbb{E}[X] = 1/\lambda$.

→ Note: the exponential distribution is memoryless.

Poisson Processes

A Poisson process $N(t)$ counts the number of arrivals with exponentially distributed inter-arrival times.

$$S_n = \sum_{i=1}^n X_i \quad X_i \sim \exp(\lambda) \quad (2.1)$$

$\forall n, t$, we have the relation $\{S_n \leq t\} = \{N(t) \geq n\}$, where S_n is a random variable telling at which time the n -th occurrence appears.

→ Note: a Poisson process is memoryless: $P[Z_1 > x] = e^{-\lambda x}$, with Z_1 be the duration of the time interval from t until the first arrival after t .

For a Poisson process of rate λ , and any given $t > 0$, the length of the interval from t until the first arrival after t is an exponentially distributed random variable. This random variable is independent of both $N(t)$ and of the $N(t)$ arrival epochs before time t . It is also independent of $N(\tau)$, $\forall \tau \leq t$.

Let us consider the process after Z_1 , Z_m , the time until the m -th arrival after time t . It is independent of $N(t)$ and of the entire previous history of the process.

Let us denote $\tilde{N}(t, t') = N(t') - N(t)$.

- Stationary increments property: It has the same distribution as $N(t' - t)$, $\forall t' \geq t$ (stationary increments property);
- Independent increments property: For any sequence of times $0 < t_1 < \dots < t_k$, the set $\{N(t_1), \tilde{N}(t_1, t_2), \dots, \tilde{N}(t_{k-1}, t_k)\}$ is a set of independent random variables.

From the memoryless property, here is another definition of a Poisson process:

- A Poisson process is a counting process that has the stationary and independent increment properties and such that

$$\begin{aligned} P[\tilde{N}(t, t + \delta) = 0] &= 1 - \lambda\delta + o(\delta) \\ P[\tilde{N}(t, t + \delta) = 1] &= \lambda\delta + o(\delta) \\ P[\tilde{N}(t, t + \delta) \geq 2] &= o(\delta) \end{aligned} \quad (2.2)$$

2.1 Distribution of $N(t)$

S_n is the sum n IID random variables and f_{S_n} is the convolution of n times f_X :

$$f_{S_n}(t) = \frac{\lambda^n t^n e^{-\lambda t}}{(n-1)!} \quad (2.3)$$

From this,

$$P[N(t) = n-1] = \frac{(\lambda t)^n e^{-\lambda t}}{(n)!} \quad (2.4)$$

and finally,

$$\mathbb{E}[N(t)] = \lambda t \quad \text{Var}[N(t)] = \lambda t \quad (2.5)$$

From equation (2.4), the Poisson process verifies the following probability conditions:

- $P[\tilde{N}(t, t+\delta) = 0] = 1 - \lambda\delta + o(\delta);$
- $P[\tilde{N}(t, t+\delta) = 1] = \lambda\delta + o(\delta);$
- $P[\tilde{N}(t, t+\delta) \geq 2] = o(\delta);$

where we use a first-order approximation of the exponential term, with $o(\delta)$ its residual. As $o(\delta)$ is negligible, we can approximate the Poisson process as a Bernoulli process.

2.1.1 Combining Poisson processes

Let $N_1(t)$ and $N_2(t)$ be two independent Poisson processes. Let the process $N(t) = N_1(t) + N_2(t)$. We can show using the three properties above that $N(t)$ is a Poisson process with rate $\lambda_1 + \lambda_2$.

2.1.2 Subdividing a Poisson process

Let $N(t)$ be a Poisson process with rate λ . We split the arrivals in 2 subprocesses $N_1(t)$ and $N_2(t)$. Each arrival of $N(t)$ is sent to $N_1(t)$ with probability p and to $N_2(t)$ with probability $(1-p)$, each split being independent from all others.

Then, the resulting processes $N_1(t)$ and $N_2(t)$ are two independent Poisson processes with respective rate $p\lambda$ and $(1-p)\lambda$.

2.1.3 Conditional arrival distribution

The density probability function when we have n Poisson processes, under the condition that $N(t) = n$, is

$$f(s_1, \dots, s_n | N(t) = n) = \frac{n!}{t^n} \quad (2.6)$$

From the previous results, we can compute that

$$P[S_1 > \tau | N(t) = n] = \left(\frac{t-\tau}{t} \right)^n \quad (2.7)$$

and the expectation is

$$E[S_1|N(t) = n] = \frac{t}{n+1} \quad (2.8)$$

And from this, we derive that

$$P[X_i > \tau|N(t) = n] = \left(\frac{t-\tau}{t}\right)^n \quad (2.9)$$

with expectation

$$E[X_i] = \frac{t}{n+1} \quad (2.10)$$

And thus the density function is

$$f_{S_i}(x|N(t) = n) = \frac{x^{i-1}(t-x)^{n-i}n!}{t^n(n-i)!(i-1)!} \quad (2.11)$$

2.2 Non-homogenous Poisson processes

A non-homogenous Poisson process $N(t)$ is a counting process with increments that are independent but not stationary, with

- $P[\tilde{N}(t, t+\delta) = 0] = 1 - \lambda(t)\delta + o(\delta);$
- $P[\tilde{N}(t, t+\delta) = 1] = \lambda(t)\delta + o(\delta);$
- $P[\tilde{N}(t, t+\delta) \geq 2] = o(\delta);$

where $\tilde{N}(t, t+\delta) = N(t+\delta) - N(t)$. The time-varying arrival rate $\lambda(t)$ is assumed to be continuous and strictly positive.

2.3 Bernoulli process approximation

We can approximate the non-homogenous Poisson process with a Bernoulli process where the time is partitioned into increments of lengths inversely proportional to $\lambda(t)$ (i.e. using a nonlinear time scale).

- $P[\tilde{N}(t, t+\epsilon/\lambda(t)) = 0] = 1 - \epsilon + o(\epsilon);$
- $P[\tilde{N}(t, t+\epsilon/\lambda(t)) = 1] = \epsilon + o(\epsilon);$
- $P[\tilde{N}(t, t+\epsilon/\lambda(t)) \geq 2] = o(\epsilon);$

Letting ϵ tend to zero, we obtain

$$P[N(t) = n] = \frac{(m(t))^n e^{-m(t)}}{n!} \quad P[\tilde{N}(t, t') = n] = \frac{(m(t, t'))^n e^{-m(t, t')}}{n!} \quad (2.12)$$

with

$$m(t) = \int_0^t \lambda(\tau) d\tau \quad m(t, t') = \int_t^{t'} \lambda(\tau) d\tau \quad (2.13)$$

2.4 Classification of queueing systems

- We note $A/B/k$ where A is the type of distribution for the arrival process, B for the service time and k the number of servers.

We suppose that the arrivals wait in a single queue. Commonly used letters are

- M: exponential distribution (for A) or Poisson process (for B);
- D: deterministic time intervals;
- E: Erlang distribution;
- G: general distribution.

Renewal Processes

A renewal process is a counting process with IID interarrival intervals. We note X_i the interval between arrivals, $\bar{X} = \mathbb{E}[X]$ is supposed to be finite with probability $P[X_i] > 0 = 1^1$, σ can be finite, and we denote $S_n = \sum_{i=1}^n X_i$ the time of the n -th arrival.

3.1 Strong law of large numbers

Let $\{N(t); t \geq 0\}$ be a renewal process, then

$$\lim_{t \rightarrow \infty} N(t) = \infty \quad \lim_{t \rightarrow \infty} \mathbb{E}[N(t)] = \infty \quad (3.1)$$

This implies that

$$\lim_{t \rightarrow \infty} \frac{N(t)}{t} = \frac{1}{\bar{X}} \text{ with probability } 1 \quad (3.2)$$

3.2 Central limit theorem

If the interarrival intervals of the renewal process $N(t)$ have a finite standard deviation, then from the CLT for IID random variables, we have

$$\lim_{t \rightarrow \infty} P \left[\frac{S_n - n\bar{X}}{\sqrt{n}\sigma} \leq \alpha \right] = \Phi(\alpha) \quad (3.3)$$

What is $\Phi(\alpha)$?

and

$$\lim_{t \rightarrow \infty} P \left[\frac{N(t) - t/\bar{X}}{\sigma\bar{X}^{-3/2}\sqrt{t}} < \alpha \right] = \Phi(\alpha) \quad (3.4)$$

→ Note: The reliability of the observed mean of successive results that are supposed to be IID depends a lot on the rule used to decide when we stop repeating the experiment.

¹A probability of 1 means that the opposite can happen, but is so rare that the probability is 0.

3.3 Stopping time

Let N be the rv corresponding to the total number of experiments observed. Let I_n be a series of rv being the indicator function of $\{N \geq n\}$:

$$I_n = \begin{cases} 1 & \text{if the } n\text{-th experiment is observed} \\ 0 & \text{otherwise} \end{cases} \quad (3.5)$$

N is a stopping time if I_n depends only on X_1, \dots, X_{n-1} . This means that stopping at 3pm, for example, is not a stopping time, because it can depend on X_n , depending if the n -th arrival is before or after 3pm.

3.3.1 Wald's inequality

Let N be a stopping time for $\{X_n; n \geq 1\}$. Then, $\mathbb{E}[S_N] = \mathbb{E}[N]\bar{X}$.

3.4 Blackwell's renewal theorem

3.4.1 Arithmetic distribution

If interarrival intervals can only have a length that is a multiple of some real number d , the interarrival distribution will be called an arithmetic distribution, and d the span of the distribution.

3.4.2 Blackwell's inequality

If the interarrival distribution of a renewal process $N(t)$ is not arithmetic, then

$$\lim_{t \rightarrow \infty} (m(t + \delta) - m(t)) = \frac{\delta}{\bar{X}} \quad \forall \delta \quad (3.6)$$

If the interarrival distribution is arithmetic with span d , then

$$\lim_{t \rightarrow \infty} (m(t + nd) - m(t)) = \frac{nd}{\bar{X}} \quad \forall n \geq 1 \quad (3.7)$$

3.4.3 Relationship with a Poisson process

The sum of many renewal processes tends to a Poisson process: for a non-arithmetic renewal process with $P[X_i = 0] = 0$, we have

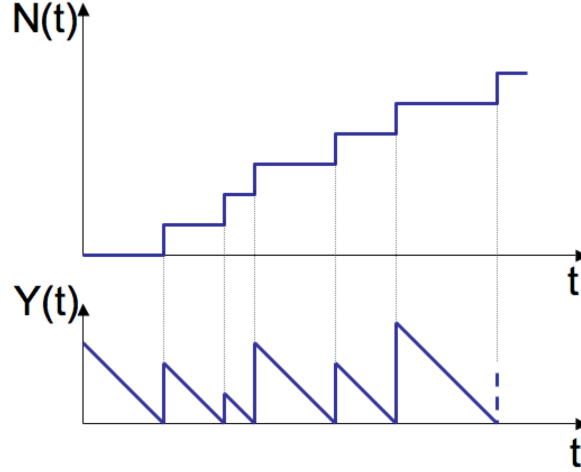
$$\begin{aligned} \lim_{t \rightarrow \infty} P[N(t + \delta) - N(t) = 0] &= 1 - \delta/\bar{X} + o(\delta) \\ \lim_{t \rightarrow \infty} P[N(t + \delta) - N(t) = 1] &= \delta/\bar{X} + o(\delta) \\ \lim_{t \rightarrow \infty} P[N(t + \delta) - N(t) \geq 2] &= o(\delta) \end{aligned} \quad (3.8)$$

The increments are asymptotically stationary, but not independent. | sectionRenewal reward process Along to the renewal process $N(t)$, we can add a reward function $R(t)$.

It models the rate at which the process is accumulating a reward or cost. It can however only depend on the current renewal but not the previous ones.

Let $Y(t)$ be the residual life at time t for the current renewal:

$$R(t) = Y(t) = S_{N(t)+1} - t \quad (3.9)$$

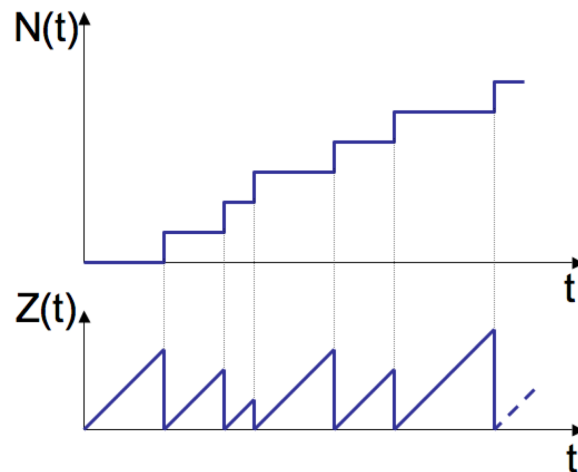


The time average residual life is $\frac{1}{t} \int_0^t Y(\tau) d\tau$.
From the definition of $Y(t)$, we can calculate that

$$\lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t Y(\tau) d\tau = \frac{\mathbb{E}[X^2]}{2\mathbb{E}[X]} = \frac{1}{2}\mathbb{E}[X] + \frac{\text{Var}(X)}{\mathbb{E}[X]} > \frac{1}{2}\mathbb{E}[X] \text{ with probability 1} \quad (3.10)$$

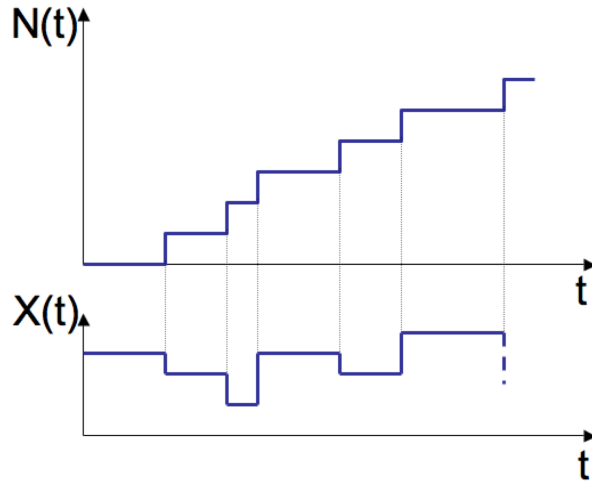
3.4.4 Time average age

Let $Z(t)$ be the age of the current renewal at time t : $R(t) = Z(t) = t - S_{N(t)}$. The time average age is $\lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t Z(\tau) d\tau = \frac{\mathbb{E}[X^2]}{2\mathbb{E}[X]}$.



3.4.5 Time average duration

Let $X(t)$ be the duration of the renewal containing time t : $R(t) = X(t) = S_{N(t)+1} - S_{N(t)}$. The time average duration is $\lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t X(\tau) d\tau = \frac{\mathbb{E}[X^2]}{\mathbb{E}[X]}$.



General renewal reward functions Let $R(t)$ be a reward function for a renewal process with expected inter-renewal times $\bar{X} < \infty$, then with probability 1,

$$\lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t R(\tau) d\tau = \frac{\mathbb{E}[R_n]}{\mathbb{E}[X]} \quad (3.11)$$

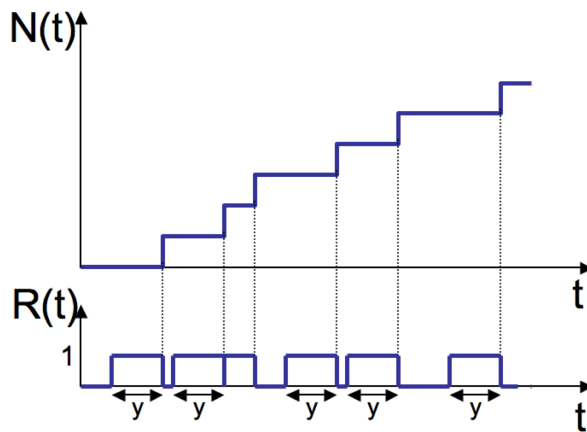
where R_n is defined as

$$R_n = \int_{S_n}^{S_{n+1}} R(\tau) d\tau \quad (3.12)$$

3.4.6 Distribution of residual life

We are interested in the fraction of time that $Y(t) \leq y$:

$$R(t) = I\{Y(t) \leq y\} \quad R_n = \min\{y, X_n\} \quad (3.13)$$



And we can calculate that

$$\begin{aligned} \mathbb{E}[R_n] &= \int_0^y P[X > x] dx \\ F_Y(y) &= \frac{1}{\mathbb{E}[X]} \int_0^y P[X > x] dx \end{aligned} \quad (3.14)$$

3.4.7 Key theorem

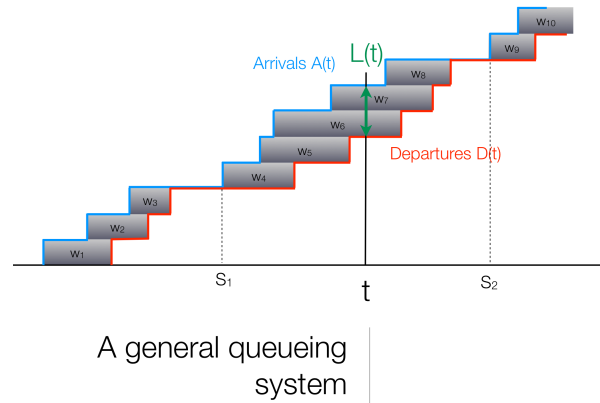
Let $N(t)$ be a non-arithmetic renewal process, let $R(z, x) \geq 0$ be such that $r(z) = \int_{x=z}^{\infty} R(z, x) dF_X(x)$ is directly Riemann integrable. Then,

$$\lim_{t \rightarrow \infty} \mathbb{E}[R(t)] = \frac{\mathbb{E}[R_n]}{\bar{X}} \quad (3.15)$$

3.5 Little's Law

Let a queueing system be such that

- $A(t)$ is the number of arrivals between 0 and t ;
- $D(t)$ is the number of departures between 0 and t ;
- $L(t) = A(t) - D(t)$ is the number of customers in the system at time t ;
- w_i the time the i^{th} customer spends in the system;
- $N(t)$ is the renewal process counting the number of busy periods of the system (each time a customer arrives when the system is empty).



Let us use $L(t)$ as a reward function for the renewal process $N(t)$. This implies

$$\begin{aligned} \sum_{n=1}^{N(t)} R_n &\leq \int_0^t L(\tau) d\tau \leq \sum_{i=1}^{A(t)} w_i \leq \sum_{n=1}^{N(t)+1} R_n \\ \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t L(\tau) d\tau &= \frac{\mathbb{E}[R_n]}{\mathbb{E}[X]} \end{aligned} \quad (3.16)$$

Putting all this together, we can show that $\bar{L} = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t L(\tau) d\tau = \bar{W}\lambda$.

3.5.1 M/G/1 queue

Let $R(t)$ be the remaining time for the customer being served. Let $U(t)$ be the time an arrival at time t would have to wait before being served. Let $L_q(t)$ be the number of

customers in queue at time t , independent of the Z_i . We define

$$U(t) = \sum_{i=1}^{L_q(t)} Z_i + R(t) \implies \mathbb{E}[U(t)] = \mathbb{E}[L_q(t)]\mathbb{E}[Z] + \mathbb{E}[R(t)] \quad (3.17)$$

We can show that

$$\int_0^{S_N(t)} R(\tau) d\tau \leq \int_0^{S_N(t)+1} R(\tau) d\tau \quad (3.18)$$

And from Little's Law,

$$\lim_{t \rightarrow \infty} \mathbb{E}[L_q(t)] = \lambda \bar{W}_q \implies \lim_{t \rightarrow \infty} \mathbb{E}[U(t)] = \lambda \bar{W}_q \mathbb{E}[Z] + \lambda \frac{\mathbb{E}[Z^2]}{2} \quad (3.19)$$

Poisson arrival process implies that arrivals occur with identical probability at any moment, this implies independence with $U(t)$. Hence $\mathbb{E}[W_q(t)] = \mathbb{E}[U(t)]$. Hence $\bar{W}_q = \lambda \bar{W}_q \mathbb{E}[Z] + \lambda \frac{\mathbb{E}[Z^2]}{2}$. And we can isolate \bar{W}_q :

$$\bar{W}_q = \frac{\lambda(\mathbb{E}[Z]^2 + \sigma^2)}{2(1 - \lambda\mathbb{E}[Z])} \quad (3.20)$$

And we remember $\bar{W} = \bar{W}_q + \mathbb{E}[Z]$.

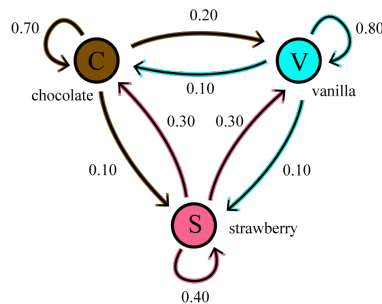
Finite State Markov Chains

4.1 Definitions

A Markov chain is a stochastic process with fixed intervals $\{X_n, n \geq 0\}$ such that each random variable $X_n, n \geq 1$ depends on the past only through the most recent random variable X_{n-1} :

$$P[X_n = j | X_{n-1} = i, X_{n-2} = k, \dots, X_0 = m] = P[X_n = j | X_{n-1} = i] = P_{ij} \quad (4.1)$$

The rv X_n is called the state of the Markov chain, and the set of possible sample values for the states lie in a countable set. A Markov chain can be represented under a graph or matrix form:



$$P = \begin{pmatrix} 0.70 & 0.20 & 0.10 \\ 0.10 & 0.80 & 0.10 \\ 0.30 & 0.30 & 0.40 \end{pmatrix} \quad (4.2)$$

- We say that a state j is accessible from i ($i \rightarrow j$) if there exists a xalk in the graph from i to j : $i \rightarrow j$ iff $P_{ij}^n = P[X_n = j | X_0 = i] > 0$ for some n .
- Two distinct states i, j communicate ($i \leftrightarrow j$) if i is accessible from j and vice versa.
- A class C of states is a non-empty set of states such that for each $i \in C$, each state $j \neq i$ satisfies $j \in C$ if $i \leftrightarrow j$ and $j \notin C$ if $i \not\leftrightarrow j$.
- A state i is recurrent if it is accessible from all states that are accessible from i . A transient state is a state that is not recurrent.

→ Note: all states of a same class are of the same type.

- A finite-state Markov chain has at least one recurrent class.
- The period of a state i , denoted $d(i)$, is the greatest common divisor of all n such that $P_{ii}^n > 0$. A state is aperiodic if $d(i) = 1$.

→ Note: All states of a class have the same periodicity.

- If a class has a period $d > 1$, then there exists a partition $\{C_i\}_{i=1}^d$ of the states of the class such that all the transitions from a state of C_n go to a state of class C_{n+1} and all transitions from C_d go to a state of C_1 , i.e. we make a cycle of subclasses.
- A class is called ergodic if it is aperiodic and recurrent.
- A matrix is stochastic iff it is square, non negative and each row sums to 1, i.e. $P\mathbb{1}_n = \mathbb{1}_n$.

4.2 Transition probabilities

We can calculate that

$$P[X_{n+2} = j | X_n = i] = P_{ij}^2 = \sum_{k=1}^J P_{ik} P_{kj} \implies P^2 = P \cdot P \implies P^n = P \cdot \dots \cdot P \quad (4.3)$$

More generally, $P_{ij}^{n+m} = \sum_{k=1}^J P_{ik}^n P_{kj}^m$.

Because of ergodicity, all rows converge to the same value, and we store those values in a row vector called π . Then, $\pi = \pi P$ and the sum of all values of π is 1.

From this, we induce that π is a left eigenvector of P for the eigenvalue 1, and the number of linearly independent solutions corresponds to the multiplicity of the eigenvalue 1. There will be one independent solution for each recurrent class of P . Moreover, if P is ergodic, then $\lim_{n \rightarrow \infty} P^n = \mathbb{1}_m \pi$, else π will be the average over the different subclasses.

4.3 Markov chains with rewards

Let r_i be the reward associated with state i . In the case where the reward is on the edge from node i to j and not the states, the reward is r_{ij} and we have $r_i = \sum_j P_{ij} r_{ij}$. For an ergodic chain, we observe that the average reward per period will be

$$g = \sum_i r_i \pi_i \quad (4.4)$$

where π_i is the steady state probability.

4.3.1 Expected reward over multiple transitions

Let X_m be the state at time m and let $R_m = R(X_m)$ be the reward at time m . Under the condition $X_m = i$ (starting point), the aggregate expected reward $v_i(n)$ over n periods from X_m to X_{m+n-1} is

$$v_i(n) = \mathbb{E}[R(X_m) + \dots + R(X_{m+n-1}) | X_m = i] = r_i + \sum_j P_{ij} r_j + \dots + \sum_j P_{ij}^{n-1} r_j \quad (4.5)$$

And in vector notation

$$v(n) = r + [P]r + \dots + [P^{n-1}]r = \sum_{h=0}^{n-1} [P^h]r \quad (4.6)$$

4.3.2 Relative gain vector

Assuming the Markov chain is an ergodic unichain, i.e. it has a single ergodic class with possibly some transient classes, we know that $\lim_{n \rightarrow \infty} [P^n] = \mathbb{1}_m \pi$ and thus

$$\lim_{n \rightarrow \infty} [P^n]r = \mathbb{1}_m \pi r = g \mathbb{1}_m \quad (4.7)$$

This means that the expected reward per period converges to g . From this, we can evaluate the transient effect and define the relative-gain vector, denoted by w :

$$w = \lim_{n \rightarrow \infty} (v(n) - n g \mathbb{1}_m) = \lim_{n \rightarrow \infty} \sum_{h=0}^{n-1} [P^h - \mathbb{1}_m \pi]r \quad (4.8)$$

It can also be computed by solving the following equations instead of calculating the limit:

$$w + g \mathbb{1}_m = [P]w + r \quad \pi w = 0 \quad (4.9)$$

The second equation means that the sum (weighted by the probabilities) of all gains is 0.

4.4 Markov decision processes

Suppose that in each state i , we can choose between K_i different possibilities with rewards $r_i^{(1)}, \dots, r_i^{(K_i)}$. This means that we choose between different Markov chains that have the same states but not necessarily the same edges. We want to find the optimal (stationary or dynamic) policy for this problem.

4.4.1 Dynamic programming algorithm for the dynamic optimal policy

This section aims to find a dynamic programming algorithm for the dynamic optimal policy. We assume here that for any policy A , the resulting Markov chain with matrix $[P^A]$ is an ergodic unichain.

Let $v(0)$ be the final reward vector. Through a recurrence method, we can show that

$$v_i^*(n) = \max_k \{r_i^{(k)} + \sum_j P_{ij}^{(k)} v_j^*(n-1)\} \quad (4.10)$$

or in vector form:

$$v^*(n) = \max_A \{r^A + [P^A]v^*(n-1)\} \quad (4.11)$$

for a policy A , i.e. a decision k_i for each state i . From equations (4.9), we know that

$$w^A = r^A - g^A \mathbb{1}_m + [P^A]w^A \quad (4.12)$$

and thus

$$v^A(n) = n g^A \mathbb{1}_m + w^A + [P^A]^n (v(0) - w^A) \quad (4.13)$$

If $v(0) = w^B$ for a policy B such that

$$r^B + [P^B]w^B \geq r^A + [P^A]w^A \quad \forall A \quad (4.14)$$

then B is the dynamic optimal policy for each time period, and

$$v^*(n) = w^B + ng^B \mathbf{1}_m \quad (4.15)$$

This relation is an equivalence.

4.4.2 Policy improvement algorithm

Algorithm 1 Policy improvement algorithm

- 1: **Step 1:** Choose an arbitrary policy B ;
 - 2: **while** $\exists A : r^B + [P^B]w^B \stackrel{\leq}{\neq} r^A + [P^A]w^B$ **do**
 - 3: **Step 2:** Compute w^B ;
 - 4: **Step 3:** Find A such that $r^A + [P^A]w^B \stackrel{\neq}{\geq} r^B + [P^B]w^B$;
 - 5: **Step 4:** $B \leftarrow A$;
 - 6: **end while**
-

Theorem 4.1. Assuming for any policy A that the Markov chain $[P^A]$ is an ergodic unichain, if B is an optimal stationary policy, then

$$\lim_{n \rightarrow \infty} v^*(n) - ng^B \mathbf{1}_m = w^B + (\beta - \pi^B w^B) \mathbf{1}_m \quad (4.16)$$

What is β ?

4.5 Dynamic programming algorithm for the stationary optimal policy

Algorithm 2 Dynamic programming algorithm for the stationary optimal policy

- 1: **Step 1:** Fix an arbitrary vector $v(0)$;
- 2: **while** $l < u$ **do**
- 3: **Step 2:** Compute

$$l = \min_i [v_i^*(n) - v_i^*(n-1)] \quad u = \max_i [v_i^*(n) - v_i^*(n-1)] \quad (4.17)$$

- 4: **end while**
 - 5: $l = u = g^A$ and A is the optimal stationary policy.
-

Markov Decision Processes and Reinforcement Learning

5.1 MDP and Policies

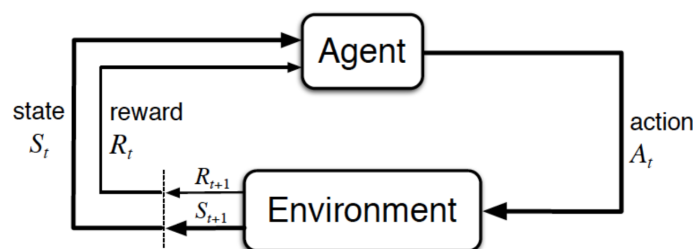
5.1.1 Markov Decision Process

A Markov Decision Process models a sequential decision-making process under uncertainty, where moving to the next stage only depends on the current action-state pair.

Definition 5.1. A Markov Decision Process (MDP) is defined by

- a set of system states S ;
- a set of actions \mathcal{A} ;
- a set of rewards R ;
- probabilities $p(s', r|s, a)$ of getting a reward r and moving to state s' if action a is taken in state s .

The MDP is finite if the three sets are finite.



From the definition, we can derive 3 quantities:

- Transition probabilities: $p(s|s, a) = \sum_r p(s', r|s, a)$;
- Reward probabilities: $p(r|s, a) = \sum_{s'} p(s', r|s, a)$;
- Expected reward knowing s, a : $r(s, a) = \mathbb{E}[R|s, a] = \sum_r r \sum_{s'} p(s', r|s, a)$.

5.1.2 Policies

A policy defines the decision making in each state;

Definition 5.2. A policy is a mapping $\pi : s \in S \rightarrow \pi(a|s)$, where $\pi(a|s)$ represents the probability of taking action a in state s .

5.1.3 Policy values and state-action values

The policy value $v_\pi(s)$ at a given state s corresponds to the expected rewards collected over time by applying policy π starting from state s .

$$v_\pi(s) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right] \quad (5.1)$$

where $\gamma \in [0, 1[$ is a discounted factor. The state-action value function is the same idea, but has a dependance on the action we intend to take.

$$q_\pi(s, a) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right] \quad (5.2)$$

We can show that

$$v_\pi(s) = \mathbb{E}_{A \sim \pi(a|s)} [q_\pi(s, A)] \quad (5.3)$$

Those definitions induce the Bellman equations, a linear system in $v_\pi(s)$:

$$\forall s \in S, \quad v_\pi(s) = \sum_a \pi(a|s) \sum_{s'|r} p(s', r | s, a) [r + \gamma v_\pi(s')] \quad (5.4)$$

5.1.4 Policy evaluation

Given the policy π and the MDP probabilities p , we can rewrite the Bellman equations as

$$V = R + \gamma P V \iff (I - \gamma P) V = R \quad (5.5)$$

where $I - \gamma P$ is invertible as $\|P\|_\infty \leq 1$ and $\gamma < 1$.

Another way to solve this is by an iterative process: introducing the linear operator $L : V \rightarrow \gamma P V + R$, we want to find V such that $V \approx L(V)$. This approach is guaranteed to converge since L is γ -Lipschitz (affine application). Defining $V^* = \lim_{t \rightarrow \infty} L^t(V_0)$, with V_0 the first iterate, we can show that

$$\|V^* - V_{n+1}\|_\infty \leq \frac{\gamma}{1 - \gamma} \|V_{n+1} - V_n\|_\infty \quad (5.6)$$

Algorithm 3 Policy Evaluation Algorithm

```
1: Input:  $\pi$  the policy to be evaluated, and  $\theta$  the guaranteed accuracy of estimation
2: Initialization:  $V(s)$  the arbitrary initial value for all  $s$ 
3: while  $\Delta \geq \theta$  do
4:    $\Delta = 0$ 
5:   for each state  $s$  do
6:      $v = V(s)$ 
7:      $V(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$ 
8:      $\Delta = \max(\Delta, |v - V(s)|)$ 
9:   end for
10: end while
```

5.2 Optimizing Policies

5.2.1 Policy Improvement Theorem

Theorem 5.3.

$$[\forall s \in S, q_\pi(s, \pi^{new}(s)) \geq q_\pi(s, \pi(s))] \implies v_{\pi^{new}}(s) \geq v_\pi(s) \quad (5.7)$$

A strict inequality on the left implies a strict one on the right too. This theorem means that a new policy will be at least as good as a given policy π if changing any action of π by the corresponding action of π^{new} yields a better total gain at the end.

5.2.2 Bellman optimality conditions

Definition 5.4. A policy π^* is optimal if for any state $s \in S$ and any other policy π , $v_{\pi^*}(s) \geq v_\pi(s)$.

Theorem 5.5. A policy π is optimal iff for any state action pair (s, a) with a positive probability to be selected by the policy, i.e. $\pi(a|s) > 0$, we have

$$a \in \arg \max_{a' \in A} q_\pi(s, a') \quad (5.8)$$

Meaning that any action with a nonzero probability to be taken maximizes the gain of the state at which it is taken.

It can be shown that this implies that any finite MDP admits an optimal policy which is deterministic.

This theorem yields the two following equations for optimal policies:

$$\begin{aligned} v_*(s) &= \max_a q_{\pi^*}(s, a) = \max_{a \in A(s)} \sum_{s',r} p(s',r|s,a)(r + \gamma v_*(s')) \\ q^*(s, a) &= \sum_{s',r} p(s',r|s,a)(r + \gamma \max_{a'} q^*(s', a')) \end{aligned} \quad (5.9)$$

5.3 Value Iteration Algorithm

In the same idea as we did in the evaluation section, let us define the operator $\Phi : V \rightarrow \Phi(V)$ such that

$$\Phi(V)(s) := \max_{a \in A(s)} \sum_{s', r} p(s', r | s, a) (r + \gamma V(s')) \quad (5.10)$$

And thus the Bellman optimality equation ((5.9)) is equivalent to $V = \Phi(V)$, which can be solved iteratively from the starting point V_0 . This method is guaranteed because Φ is a contracting operator (\sim L-Lipschitz function).

Algorithm 4 Value Iteration Algorithm

```
1: Input: the guaranteed accuracy of estimation  $\theta$ ;  
2: Initialization:  $V(s) :=$  an arbitrary initial value for all  $s$ ;  
3:  $\Delta \geq \theta$ ;  
4: while  $\Delta \geq \theta$  do  
5:    $\Delta = 0$ ;  
6:   for each state  $s$  do  
7:      $v = V(s)$   
8:      $V(s) = \max_a \sum_{s', r} p(s', r | s, a) (r + \gamma V(s'))$ ;  
9:      $\Delta = \max(\Delta, |v - V(s)|)$ ;  
10:   end for  
11: end while
```
