

Course structure:

① Probability methods

- 1.1. Probability generating functions
- 1.2. Inequalities and limit theorems
- 1.3. Stochastic processes

② Fourier & related methods

- 2.1. Fourier representations
- 2.2. Discrete Fourier Methods
- 2.3. Wavelets

1. Probability methods

(1.1) Probability generating functions (PGFs)

Given a random variable $X \in N$, s.t. $P(X=k) = p_k$, define $(k \in N)$

the probability generating function as:

$$G_X(z) = \sum_{k=0}^{\infty} p_k z^k$$

for all such z , for which the sum converges.

This function "encodes" all p_k 's and so defines the RV completely.

Properties:

1.) $G_X(0) = p_0$

2.) $G_X(1) = \sum_{k=0}^{\infty} p_k = 1$

3.) $G_X(z) = \sum_{k=0}^{\infty} p_k z^k = \sum_{k=0}^{\infty} P(X=k) \cdot z^k \Big|_{X=k} = E(z^X)$.

4.) PGF is defined for all $|z| \leq 1$,

since $\sum_{k=0}^{\infty} |p_k z^k| \leq \sum_{k=0}^{\infty} p_k = 1$ and abs. convergence \Rightarrow convergence

4.) PGF characterizes distribution of RV, namely

$$G_X(z) = G_Y(z) \iff \forall k \in \{0, 1, 2, \dots\}. P(X=k) = P(Y=k).$$

Example: For binomial distribution

$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}, \text{ so PGF:}$$

$$G_X(z) = \sum_{k=0}^n p^k (1-p)^{n-k} \binom{n}{k} \cdot z^k = \sum_{k=0}^n \binom{n}{k} (pz)^k (1-p)^{n-k} = (pz + 1 - p)^n.$$

Uses of PGFs:

- 1.) Working with sums of IID RVs
(such as sum of a pair of dice)
- 2.) Finding moments (mean, variance etc) using differentiation

Theorem

$$\frac{d^x}{dz^x} G_X(z) = \mathbb{E}[X(X-1)\dots(X-x+1)]$$

(proven by induction)

It follows that:

$$1) \mathbb{E}(X) = G'_X(1).$$

$$\begin{aligned} 2) \text{Var}(X) &= \mathbb{E}[(X - \mathbb{E}(X))^2] = \mathbb{E}[X^2 - 2X\mathbb{E}(X) + \mathbb{E}^2(X)] = \\ &= \mathbb{E}(X^2) - \mathbb{E}^2(X). = \mathbb{E}[X(X-1)] + \mathbb{E}(X) - \mathbb{E}^2(X) = \\ &= \underline{G''(1) + G'(1)} - \underline{(G'(1))^2}. \end{aligned}$$

Theorem

$$G_{X+Y}(z) = G_X(z) \cdot G_Y(z) \quad \text{for independent } X \text{ and } Y$$

$$\text{Proof: } G_X(z) \cdot G_Y(z) = \mathbb{E}(z^X) \cdot \mathbb{E}(z^Y) = \mathbb{E}(z^{X+Y}) = G_{X+Y}(z).$$

Hence, ~~$G_{X+Y}(z)$~~

$$\mathbb{E}(X+Y) = G'_X(1) \cdot G_Y(1) + G_X(1) \cdot G'_Y(1) = \mathbb{E}(X) + \mathbb{E}(Y).$$

$$\text{Var}(X+Y) = \dots$$

(1.2)

Limits and inequalities

Topic structure:

- 1.) Moment generating functions
- 2.) Probabilistic inequalities
- 3.) Limit theorems

1.2.1 Moment generating functions

Def For a RV X define the moment generating function as:

$$M_x(t) = \mathbb{E}(e^{tX})$$

for those $t \in \mathbb{R}$ for which the expectation exists.

Properties: 1.) $M_x(0) = 1$.

2.) If X has MGF $M_x(t)$ and $y = ax + b$, then

$$M_y(t) = \mathbb{E}(e^{t(ax+b)}) = e^{bt} \cdot M_x(at).$$

3.) If X and Y are independent RVs then

$$M_{X+Y}(t) = \mathbb{E}(e^{t(X+Y)}) = \mathbb{E}(e^{tX} \cdot e^{tY}) = \mathbb{E}(e^{tX}) \cdot \mathbb{E}(e^{tY}) = M_X(t) \cdot M_Y(t).$$

4.) Since $M_X(t) = \mathbb{E}(e^{tX}) = \mathbb{E}(1 + tX + \frac{(tX)^2}{2!} + \dots) = 1 + t \cdot \mathbb{E}(X) + \frac{t^2}{2!} \mathbb{E}(X^2) + \dots$,
 $\left(\frac{d^n}{dt^n} M_X(t) \right) \Big|_{t=0} = \mathbb{E}(X^n)$ (n -th moment of X).

5.) If X is discrete RV $\in \{0, 1, 2, \dots\}$ with PGF $G_X(t)$, then

$$M_X(t) = \mathbb{E}(e^{tX}) = G_X(e^t).$$

6.) Uniqueness: same MGF \Leftrightarrow same distribution function

7.) Continuity: $[\lim_{n \rightarrow \infty} F_n(x) = F(x)] \Leftrightarrow [\lim_{n \rightarrow \infty} M_n(x) = M(x)]$

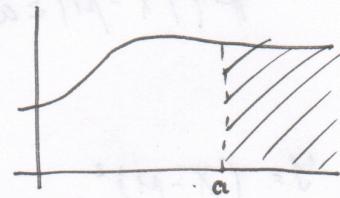
where $F_n(x)$ - probability distrib. functions

$M_n(x)$ - corresponding MGFs

1.2.2 Probabilistic inequalities

It is generally useful to have upper bounds our probabilities of the form:

$$P(X \geq a)$$



We'll consider 3 such inequalities:

- Markov's inequality
- Chebychev's inequality
- Chernoff's inequality

useful concept:

$$I(A) = \begin{cases} 1, & \text{if } A \text{ happens} \\ 0, & \text{otherwise} \end{cases} \quad (\text{indicator})$$

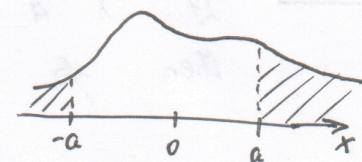
$$\text{Note: } E(I(A)) = P(A) \cdot 1 + P(\bar{A}) \cdot 0 = P(A)$$

1.) Markov's inequality

Theorem If $E(X) < \infty$, then for any $a > 0$

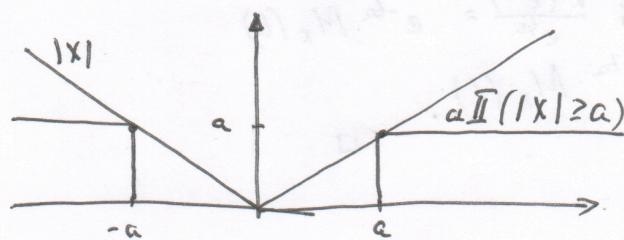
$$P(|X| > a) \leq \frac{E(|X|)}{a}$$

Proof



$$I(|X| \geq a) = \begin{cases} 1, & \text{if } |X| \geq a \\ 0, & \text{if } -a < X < a \end{cases}$$

Consider the graphs of $|X|$ and $aI(|X| \geq a)$.



From the drawing, it follows that

$$|X| \geq a \cdot I(|X| \geq a)$$

Taking expectations of both sides:

$$E(|X|) \geq a \cdot P(|X| \geq a)$$

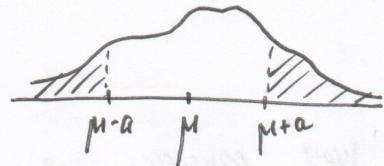
$$P(|X| \geq a) \leq \frac{E(|X|)}{a}$$

a.) Chebychev's inequality

Theorem

If X is a RV with mean μ and finite variance σ^2 , then for all $a > 0$,

$$P(|X - \mu| \geq a) \leq \frac{\sigma^2}{a^2}$$



Proof Let $Y = (X - \mu)^2$

$$\text{Then } E(Y) = E((X - \mu)^2) = \text{Var}(X) = \sigma^2$$

$$\text{By Markov: } P(|Y| \geq b) \leq \frac{E(Y)}{b}$$

$$P((X - \mu)^2 \geq b) \leq \frac{\sigma^2}{b}$$

$$\text{Now let } b = a^2, a > 0$$

$$P((X - \mu)^2 \geq a^2) = P(|X - \mu| \geq a) \leq \frac{\sigma^2}{a^2} \quad \square$$

b.) Chernoff's inequality

Theorem

If X is defined by MGF $M_x(t)$ and $a \in \mathbb{R}$, then for all $t > 0$,

$$P(X \geq a) \leq e^{-ta} \cdot M_x(t).$$

Note: we can choose $t > 0$ to make the bound as tight as possible

Proof

$$P(X \geq a) = P(e^{tX} \geq e^{ta})$$

$$\text{By Markov: } P(e^{tX} \geq e^{ta}) \leq \frac{E(e^{tX})}{e^{ta}} = e^{-ta} \cdot M_x(t)$$

$$\text{So } P(X \geq a) \leq e^{-ta} \cdot M_x(t).$$

□

1.2.3

Limit theorems
Topic structure
1) Convergence

(4)

- a) Weak Law of Large Numbers (WLLN)
- b) Central Limit Theorem (CLT)
- Application: calculation of confidence intervals

1.) Convergence

For a sequence of random variables $\{X_n\}_{n \geq 1}$, we shall define two notions of convergence to some RV X as $n \rightarrow \infty$:

a.) Convergence in distribution ($X_n \xrightarrow{D} X$)

Def: $X_n \xrightarrow{D} X \iff F_{X_n} \rightarrow F_X$ as $n \rightarrow \infty$.

b.) Convergence in probability ($X_n \xrightarrow{P} X$)

Def: $X_n \xrightarrow{P} X \iff P(|X_n - X| < \epsilon) \rightarrow 0$ as $n \rightarrow \infty$

Theorem

Convergence in probability implies convergence in distribution
(convergence in probability is a stronger statement)

Proof $X_n \xrightarrow{P} X \Rightarrow X_n \xrightarrow{D} X$

Consider the cumulative distribution function of X_n

$$F_{X_n}(a) = P(X_n \leq a)$$

Now

$$F_{X_n}(a) = P(X_n \leq a) = P(X_n \leq a \wedge X > a + \epsilon) + P(X_n \leq a \wedge X \leq a + \epsilon) \quad (\text{for an arbitrary fixed } \epsilon > 0) \quad (I)$$

If $\begin{cases} X_n \leq a \\ X > a + \epsilon \end{cases} \iff \begin{cases} -X_n \geq -a \\ X > a + \epsilon \\ X - X_n > \epsilon \end{cases}$, then $X - X_n > \epsilon \Rightarrow |X - X_n| > \epsilon$. (II)

Also: $X_n \leq a \wedge X \leq a + \epsilon \Rightarrow X \leq a + \epsilon$, so $P(X_n \leq a \wedge X \leq a + \epsilon) \leq P(X \leq a + \epsilon)$ (III)

Putting (II) and (III) into (I), we get

$$F_{X_n}(a) \leq P(|X - X_n| > \epsilon) + P(X \leq a + \epsilon) = P(|X - X_n| > \epsilon) + F_X(a + \epsilon).$$

~~Similarly, $F_{X_n}(a) \geq P(|X - X_n| < \epsilon)$~~

$$\text{Similarly, } F_{X_n}(a-\epsilon) = P(X \leq a-\epsilon \wedge X_n > a) + P(X \leq a-\epsilon \wedge X_n \leq a) \leq \\ \leq P(|X_n - X| > \epsilon) + F_X(a)$$

Using the two inequalities, we can sandwich $F_{X_n}(a)$:

$$F_X(a-\epsilon) - P(|X_n - X| > \epsilon) \leq F_{X_n} \leq P(|X_n - X| > \epsilon) + F_X(a+\epsilon)$$

since $X_n \xrightarrow{P} X$, $P(|X_n - X| > \epsilon) \rightarrow 0$ (by definition),

so

$$F_X(a-\epsilon) \leq F_{X_n}(a) \leq F_X(a+\epsilon) \quad \text{for an arbitrary small } \epsilon$$

so $F_{X_n} \rightarrow F_X$, which is the definition of convergence in dist.

Limit theorems

Idea: Given a sequence of RVs $(X_n)_{n \geq 1}$, let

$$S_n = X_1 + X_2 + X_3 + \dots + X_n \quad \text{and} \quad \bar{X}_n = \frac{S_n}{n}$$

What happens to the sample average \bar{X}_n for a large n ?

d.) The Weak Law of Large Numbers (WLLN)

Theorem If $(X_n)_{n \geq 1}$ is a sequence of IID RVs with finite mean μ , then $\bar{X}_n \xrightarrow{P} \mu$.

Proof

$$\begin{aligned} E(\bar{X}_n) &= E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \frac{1}{n} \sum_{i=1}^n \mu = \mu \\ \text{Var}(\bar{X}_n) &= \text{Var}\left(\frac{X_1 + X_2 + \dots + X_n}{n}\right) = \text{Var}\left(\frac{X_1}{n}\right) + \text{Var}\left(\frac{X_2}{n}\right) + \dots + \text{Var}\left(\frac{X_n}{n}\right) = \\ &= \underbrace{\frac{\sigma^2}{n^2} + \frac{\sigma^2}{n^2} + \dots + \frac{\sigma^2}{n^2}}_{n \text{ terms}} = \frac{\sigma^2}{n} \end{aligned}$$

By Chebyshev's inequality:

$$0 \leq P(|\bar{X}_n - \mu| > \epsilon) \leq \frac{\sigma^2}{n \cdot \epsilon^2}, \text{ for an arbitrary small } \epsilon$$

Hence, $\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| > \epsilon) = 0$, so $\bar{X}_n \xrightarrow{P} \mu$
for large n . \square

3.) Central Limit Theorem (CLT)

(5)

Theorem For a sequence of IID RVs $(X_n)_{n \geq 1}$ with mean μ and variance σ^2 , whose MGF converges in some interval $(-\alpha, \alpha)$ ($\alpha > 0$),

$$Z_n = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow{D} Z \sim N(0, 1)$$

What it means: Given that some conditions hold, the sample average (\bar{X}_n) is distributed roughly normally (in a bell curve). That is, if we compute the averages of many samples and then plot them, we'll get the familiar bell curve.

Proof Let $y_i = \frac{X_i - \mu}{\sigma}$, then

$$E(y_i) = E\left(\frac{X_i - \mu}{\sigma}\right) = \frac{E(X_i) - \mu}{\sigma} = \frac{\mu - \mu}{\sigma} = 0.$$

$$E(y_i^2) = E[(y_i - E(y_i))^2] = \text{Var}(y_i) = 1$$

Hence, the MGF:

$$M_{y_i}(t) = 1 + \underbrace{\frac{t^2}{2} + o(t^2)}_{\text{higher order terms}}$$

Now,

$$\bar{Z}_n = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} = \frac{\frac{1}{n} \sum_{i=1}^n X_i - \mu}{\sigma/\sqrt{n}} = \frac{\sum_{i=1}^n \left(\frac{X_i}{n} - \frac{\mu}{n}\right)}{\sigma/\sqrt{n}} = \frac{\sqrt{n}}{n} \sum_{i=1}^n \frac{X_i - \mu}{\sigma} = \frac{1}{\sqrt{n}} \sum_{i=1}^n y_i$$

Combining properties (2) and (3) of MGFs:

$$M_{\bar{Z}_n}(t) = (M_{y_i}\left(\frac{t}{\sqrt{n}}\right))^n = \left(1 + \frac{t^2}{2n} + o\left(\frac{t^2}{n}\right)\right)^n$$

$$\lim_{n \rightarrow \infty} M_{\bar{Z}_n}(t) = \lim_{n \rightarrow \infty} \left(1 + \frac{t^2}{2n} + o\left(\frac{t^2}{n}\right)\right)^n = e^{\frac{t^2}{2}}$$

which is exactly the MGF of $N(0, 1)$.

□

Applications of CLT: Confidence Intervals

Confidence interval - interval estimate of some unknown value, accurate with some known, parametrically controllable probability.

Theory

By CLT

$$Z_n = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \quad (\text{I})$$

is asymptotically distributed as $N(0, 1)$.

In real life, we can't know the true value of the variance σ^2 , so instead we use the so-called sample variance:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{X}_i - \bar{X}_n)^2$$

We can do it because we can show that $E(S^2) = \sigma^2$. So, instead of (I), we use

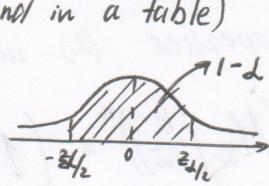
$$Z_n = \frac{\bar{X}_n - \mu}{S/\sqrt{n}}$$

which is also approximately distributed as $N(0, 1)$. (II)

Now let's define a z_α , so that $P(Z > z_\alpha) = \alpha$ for $Z \sim N(0, 1)$

(values of z_α can usually be found in a table)

Then $P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha$



From (II), we can substitute Z_n since it is approximately distributed as $N(0, 1)$:

$$P(-z_{\alpha/2} < \frac{\bar{X}_n - \mu}{S/\sqrt{n}} < z_{\alpha/2}) \approx 1 - \alpha$$

$$P(\bar{X}_n - z_{\alpha/2} \frac{S}{\sqrt{n}} < \mu < \bar{X}_n + z_{\alpha/2} \frac{S}{\sqrt{n}}) \approx 1 - \alpha$$

So the true mean μ lies within the interval $\bar{X}_n \pm z_{\alpha/2} \frac{S}{\sqrt{n}}$ with probability $1 - \alpha$. We can vary α to make interval tighter or increase coverage probability.

(1.3) Stochastic processes

(6)

Topic structure:

1.) Random walks

- a. Intro and basic concepts
- b. Returning to starting state
- c. Gambler's ruin problem

2.) Markov chains

- a. Basic definitions, concepts and theorems

~~etc.~~

- b. Hidden Markov Models (HMMs) and associated problems

1.3.1. Random Walks

a.) Introduction and basic concepts

Consider a sequence y_1, y_2, \dots of IID RVs such that:

$$P(y_i = 1) = p, P(y_i = -1) = 1-p \text{ for } p \in [0; 1].$$

(so y_i ~~can~~ only takes values +1 or -1 with probabilities p and $1-p$).

Definition A simple random walk is a sequence of RVs $\{X_n | n \in \mathbb{Z}^+\}$ defined by

$$X_n = X_0 + \sum_{i=1}^n y_i$$

where X_0 is the starting value

(so $\sum_{i=1}^n y_i$ is the sum of several ± 1 steps)

Definition A simple symmetric random walk is that where y_i is equally likely to be +1 and -1 (so $p = \frac{1}{2}$).

b.) Returning to the starting state

What is the probability of returning after exactly n steps?

$$P(X_n = X_0) = \begin{cases} 0, & n \text{ is odd} \\ \binom{n}{m} p^m (1-p)^m & \end{cases}$$

Need the same number $\binom{m}{m}$ of up- and down-steps.

What is the probability of ever returning to the starting state?

A simple random walk is called recurrent if it is certain to come back to its starting state at some time in the future. Otherwise, it is called transient.

Theorem

For a simple random walk starting at $X_0 = 0$, the probability of ever revisiting the starting state is:

$$P(X_n = 0 \text{ for some } n \in \mathbb{Z}^+) = 1 - |\mathbb{E}(y_i)| = 1 - |2p - 1|.$$

Interesting fact

$$\mathbb{E}(y_i) = 1 \cdot p - 1(1-p) = 2p - 1$$

if $p > \frac{1}{2}$, $\mathbb{E}(y_i) > 0 \rightarrow$ net drift upwards

if $p < \frac{1}{2}$, $\mathbb{E}(y_i) < 0 \rightarrow$ net drift downwards

if $p = \frac{1}{2}$, $\mathbb{E}(y_i) = 0 \rightarrow$ no net drift

Note: this implies that a simple random walk can only be ~~ever~~ recurrent if it is symmetric ($p = \frac{1}{2}$), and it is transient if $p \neq \frac{1}{2}$.

Proof:

Without loss of generality, assume that $X_0 = 0$. Define event $R_n = \{X_n = 0\}$, that is that the simple random walk returns to X_0 at time n . Now consider event:

$$\mathcal{F}_n = \{X_n = 0, X_m \neq 0 \text{ for } m \in \mathbb{Z}^+, m < n\}$$

which is that the random walk first revisits its starting state at time n .

Then:

$$P(R_n) = \sum_{m=1}^n P(R_n \cap \mathcal{F}_m)$$

but

$$P(R_n \cap \mathcal{F}_m) = P(\mathcal{F}_m) \cdot P(R_{n-m})$$

$$\text{So } P(R_n) = \sum_{m=1}^n P(\mathcal{F}_m) \cdot P(R_{n-m}), \text{ or, more briefly:}$$

$$\gamma_n = \sum_{m=1}^n f_m \cdot \gamma_{n-m} \quad (\text{I})$$

Now let's consider the probability generating functions for γ_n and f_n : (7)

$$R(z) = \sum_{n=0}^{\infty} \gamma_n z^n \quad F(z) = \sum_{n=0}^{\infty} f_n z^n \quad |z| < 1$$

$$\gamma_0 = 1 \\ (\text{we start at } X_0)$$

$$f_0 = 0$$

\downarrow
to make the series converge

Now let's derive an expression for $F(z)$ in terms of $R(z)$:

$$\begin{aligned} R(z) - 1 = \sum_{n=1}^{\infty} \gamma_n z^n &= \underbrace{\sum_{n=1}^{\infty} \sum_{m=1}^n f_m \gamma_{n-m} z^n}_{\text{shifting over the same range of } m \text{'s and } n \text{'s}} = \\ &= \underbrace{\sum_{m=1}^{\infty} \sum_{n=m}^{\infty} f_m z^m \gamma_{n-m} z^{n-m}}_{\text{Substitution } k=n-m} = \\ &= \sum_{m=1}^{\infty} f_m z^m \sum_{k=0}^{\infty} \gamma_k z^k = (F(z) - 1) \cdot R(z). \end{aligned}$$

shifting over the same range of m 's and n 's

$$\text{So } R(z) - 1 = F(z) \cdot R(z)$$

$$F(z) = 1 - R^{-1}(z). \quad (II)$$

Now let's explicitly calculate $R(z)$:

$$\begin{aligned} R(z) &= \sum_{n=0}^{\infty} \gamma_n z^n = \sum_{m=0}^{\infty} \gamma_{2m} z^{2m} = \quad (\text{since } \gamma_{2m+1} = 0 \quad \forall m \in \mathbb{Z} \geq 0) \\ &= \sum_{m=0}^{\infty} \binom{dm}{m} p^m (1-p)^m z^{2m} = \\ &= \sum_{m=0}^{\infty} \binom{2m}{m} (p(1-p)z^2)^m = \\ &= (1 - 4p(1-p)z^2)^{-\frac{1}{2}} \quad (\text{can expand it to see it gives same series as previous line}) \end{aligned}$$

By (II):

$$F = 1 - (1 - 4p(1-p)z^2)^{\frac{1}{2}}$$

Finally:

$$\begin{aligned} \# \Pr(X_n=0, \text{ for some } n \in \mathbb{Z}^+) &= \Pr(F_1 \cup F_2 \cup F_3 \dots) = f_1 + f_2 + f_3 + \dots = \\ &= \lim_{z \rightarrow 1^-} \sum_{n=1}^{\infty} f_n z^n = F(1) = 1 - ((2p-1)^2)^{\frac{1}{2}} = 1 - |2p-1| \end{aligned}$$

□

Mean return time

Suppose $p = \frac{1}{2}$, so the simple random walk is recurrent.

Let's say:

$$T = \min\{n \geq 1 | X_n = 0\} \quad \text{(time of first return)}$$

Then

$$P(T = n) = f_n$$

So the PGF for T is the same as that for F (in theorem before). Namely:

$$G_T(z) = 1 - (1 - 4p(1-p)z^2)^{\frac{1}{2}} = 1 - (1-z^2)^{\frac{1}{2}} \quad \begin{matrix} \text{since } p = \frac{1}{2}, \\ |z| < 1 \end{matrix} \quad \begin{matrix} \text{we are considering} \\ \text{the symmetric case} \end{matrix}$$

By definition of PGF, the mean of return time T :

$$E(T) = \left. \frac{d}{dz} G_T(z) \right|_{z=1} = \lim_{z \rightarrow 1^-} z (1-z^2)^{-\frac{1}{2}} = \infty$$

→ since we said $|z| < 1$.

So the symmetric simple random walk is certain to return to starting state, but the mean return time is infinite. Such random walks are called null-recurrent.

c.) The gambler's ruin problem

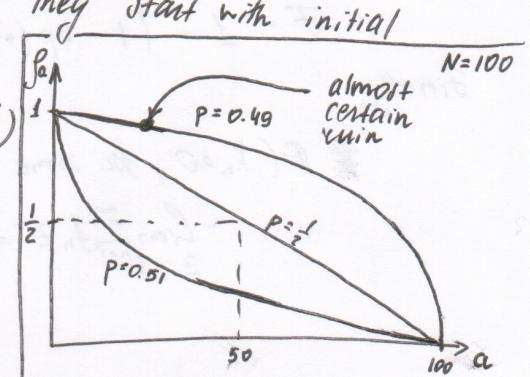
Problem: Two players (A and B) with joint capital between them £N.

Player A starts with £a ($0 \leq a \leq N$) of capital and at each time step either receives £1 from B with probability p , or gives £1 to B with probability $q = 1-p$. What is the probability of ruin (losing the whole capital) for Player A?

Theorem: The probability of ruin for A when they start with initial capital a is given by:

$$f_a = \begin{cases} \frac{\theta^a - \Theta^N}{1 - \Theta^N} & \text{if } P \neq q \text{ (i.e. } p \neq \frac{1}{2}) \\ 1 - \frac{a}{N} & \text{if } P = q \text{ (i.e. } p = \frac{1}{2}) \end{cases}$$

$$\text{where } \Theta = \frac{q}{p}.$$



Proof

$$\underline{f_a} = P(y_i=+1) \cdot f_{a+1} + P(y_i=-1) \cdot f_{a-1} = \underline{P \cdot f_{a+1} + q \cdot f_{a-1}}$$

The above is called a difference equation, and is solved by methods similar to those used for differential equations.

Try a solution of the form $f_a = \lambda^a$:

$$\lambda^a = p \cdot \lambda^{a+1} + q \cdot \lambda^{a-1}$$

$$p\lambda^2 - \lambda + q = 0.$$

$$\lambda = 1 \text{ or } \lambda = \frac{q}{p}.$$

if $p \neq q$

general solution of the form $A + B \left(\frac{q}{p}\right)^a$
Applying boundary conditions

$$f_0 = 1 \quad f_N = 0$$

We get

$$f_a = \frac{\left(\frac{q}{p}\right)^a - \left(\frac{q}{p}\right)^N}{1 - \left(\frac{q}{p}\right)^N}$$

$\left. \begin{array}{l} \text{if } p=q \\ \text{general solution of the form } C + Da \\ \text{Applying boundary conditions} \\ f_0 = 1 \text{ and } f_N = 0 \\ \text{we get} \\ f_a = 1 - \frac{q}{N}. \end{array} \right\}$

Mean duration time

How long (on average) will it take before one of the players gets ruined?
Let M_a be mean duration time for the case when A starts with f_a .
The, similarly to above:

$$M_a = 1 + P \cdot M_{a+1} + q \cdot M_{a-1},$$

Using methods similar to above it can be now shown that:

$$M_a = \begin{cases} \frac{1}{p-q} \left(N \frac{\left(\frac{q}{p}\right)^a - 1}{\left(\frac{q}{p}\right)^N - 1} - a \right) & \text{if } p \neq q \\ a(N-a) & \text{if } p=q \end{cases}$$

1.3.2. Markov Chains (MCs)

a.) Basic definitions, concepts and theorems

A Markov chain is a sequence of discrete RVs taking values in some countable state space S with the property:

$$P(X_n = a_n | X_0 = a_0, X_1 = a_1, \dots, X_{n-1} = a_{n-1}) = P(X_n = a_n | X_{n-1} = a_{n-1}).$$

→ this means that the next state is only determined by the current state and not the previous ones (MCs are memoryless)

Since S is countable, we usually map its elements to integers and "mark" the states with those.

The dynamics of a ~~the~~ Markov chain is determined by transition probabilities of the form $P(X_n = i | X_{n-1} = j)$. Since we usually want to know the probabilities of transition between any 2 states in the MC, it is convenient to define the transition matrix P , the elements of which are:

$$(P)_{ij} = p_{ij} = P(X_n = j | X_{n-1} = i)$$

where i, j are in S .

This matrix P completely defines the behaviour of the MC.

P is called a stochastic matrix (it contains probabilities), so the sum of each of its rows is 1 (total probability of going anywhere).

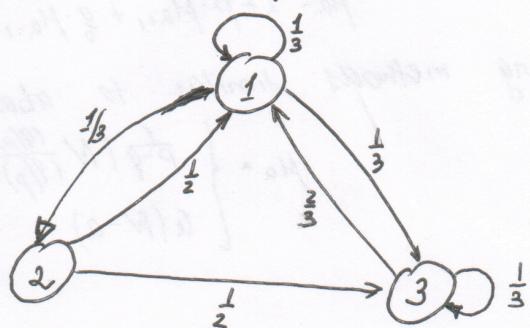
If the probabilities in the matrix do not depend on time, the MC is said to be time-homogeneous.

row - current state
column - next state

Example

$$S = \{1, 2, 3\}$$

$$P = \begin{pmatrix} 1/3 & 1/3 & 1/3 \\ 1/2 & 0 & 1/2 \\ 2/3 & 0 & 1/3 \end{pmatrix}$$



We also define the n-step transition matrix $P^{(n)} = (p_{ij}^{(n)})$ where

$$p_{ij}^{(n)} = P(X_n = j | X_0 = i)$$

so it consists of probabilities of moving from i to j in exactly n steps.

Theorem (Chapman-Kolmogorov)

For all states i, j and steps m, n :

$$P_{ij}^{(m+n)} = \underbrace{\sum_k P_{ik}^{(m)} \cdot P_{kj}^{(n)}}_{\text{matrix multiplication}}$$

Equivalently, $P_{ij}^{(m+n)} = P_{ij}^{(m)} \cdot P_{ij}^{(n)}$ $P_{ij}^{(n)} = P^n$ we can multiply matrices together to obtain the n -step transition matrix

Proof

$$\begin{aligned} P_{ij}^{(m+n)} &= P(X_{m+n} = j | X_0 = i) = \sum_k P(X_{m+n} = j, X_m = k | X_0 = i) = \\ &= \sum_k P(X_{m+n} = j | X_m = k, X_0 = i) \cdot P(X_m = k | X_0 = i) = \\ &= \sum_k P(X_{m+n} = j | X_m = k) \cdot P(X_m = k | X_0 = i) = \sum_k P_{ik}^{(n)} \cdot P_{kj}^{(m)}. \end{aligned}$$

introducing a state in the middle

Related Lemma

If $\lambda^{(n)}$ is a row vector containing probabilities of being in each state after n steps, then

$$\lambda^{(m+n)} = \lambda^{(m)} \cdot P^{(n)}$$

Proof:

$$\begin{aligned} \lambda_j^{(m+n)} &= P(X_{m+n} = j) = \sum_i P(X_{m+n} = j | X_m = i) \cdot P(X_m = i) = \\ &= \sum_i P(X_n = j | X_0 = i) \cdot P(X_m = i) = \sum_i P_{ij}^{(n)} \cdot \lambda_i^{(m)} = (\lambda^{(m)} \cdot P^{(n)})_j. \end{aligned}$$

Classification of states

Definition

- 1) j is accessible from i ($i \rightarrow j$) iff $\exists n \geq 0. P_{ij}^{(n)} > 0$
- 2) i and j communicate ($i \leftrightarrow j$) iff $(i \rightarrow j) \wedge (j \rightarrow i)$

The communicates relation (\leftrightarrow) is reflexive, symmetric and transitive, so it is an equivalence relation.

We can partition the state space into subsets in which every state communicates with all others. These subsets are called communicating classes.

Definitions

- 1.) A communicating class which, once entered, cannot be left is called closed. 
- 2.) A closed communicating class containing one ~~closed~~ state is called absorbing. 
- 3.) When the state space forms a single communicating class (can get to every state from every state), the Markov chain is called irreducible. Otherwise, it's reducible.

Recurrence and Transience of MCs

Let's define the matrix $F^{(n)} = (f_{ij}^{(n)})$ where

$$f_{ij}^{(n)} = \text{IP}(X_1 \neq j, X_2 \neq j, \dots, X_n = j | X_0 = i)$$

so $f_{ij}^{(n)}$ is the probability of visiting state j for the first time at n^{th} step, after starting from state i . Then we can say that

$$f_{ij} = \sum_{n=1}^{\infty} f_{ij}^{(n)}$$

is the probability of ever visiting j starting from i .

If $f_{ii} < 1$, the state i is transient

(meaning that we are not guaranteed to ever revisit it after visiting once)

If $f_{ii} = 1$, the state i is recurrent

(meaning that we are certain to revisit it)

Since MCs are memoryless, the probability of coming back to state i at least N times is $f_{ii} f_{ii} f_{ii} \dots f_{ii} = f_{ii}^N$ (because each return is an independent event). Then the probability of coming back to i an infinite number of times is:

$$P_{\infty} = \begin{cases} 1^{\infty} = 1, & \text{if } f_{ii} = 1 \\ f_{ii}^{\infty} = 0, & \text{if } f_{ii} < 0. \end{cases}$$

State i is transient $\Leftrightarrow \sum_{n=1}^{\infty} p_{ii}^{(n)}$ converges

State i is recurrent $\Leftrightarrow \sum_{n=1}^{\infty} p_{ii}^{(n)}$ diverges

(can be used as a criterion for recurrence/transience)

Theorem

If i and j belong to the same communicating class, then they are either both recurrent or both transient. (the solidarity property)

Proof - Slide 45 in the notes.

Mean recurrence time

Let's say

$$T_j = \min\{n \geq 1 : X_n = j\}$$

is the first time we visit j . $T_j = \infty$ if it never happens.

The mean recurrence time:

$$M_i = E(T_i | X_0 = i) = \sum_{n=1}^{\infty} n \cdot P(T_i = n | X_0 = i) = \begin{cases} \sum_{n=1}^{\infty} n \cdot f_{ii}^{(n)}, & \text{if } i \text{ is recurrent} \\ \infty \cdot \underbrace{P(T_i = \infty | X_0 = i)}_{> 0} + \dots = \infty, & \text{if } i \text{ is transient} \end{cases}$$

N.B.: $\sum_{n=1}^{\infty} n \cdot f_{ii}^{(n)}$ may not converge, so the mean recurrence time may be infinite, even if i is recurrent.

- To summarise, there are 3 possible cases. A state i can be:
- Transient - $f_{ii} < 1$, $M_i = \infty$
 - Null-recurrent - $f_{ii} = 1$, $M_i = \infty$
 - Positive recurrent - $f_{ii} = 1$, $M_i < \infty$
- $\} \text{recurrent}$

Periodicity

Let

$$d_i = \text{GCD}\{n : p_{ii}^{(n)} > 0\}$$

if $d_i = 1$, the state is aperiodic
else if $d_i > 1$, the state is periodic

It can be shown that all states in a communicating class have same period.

From here onwards, we only consider irreducible, aperiodic Markov chains (i.e. consisting of a single communicating class with aperiodic states).

Stationary distribution

Stationary (a.k.a. equilibrium) distribution is basically:

$$\pi = \lim_{n \rightarrow \infty} \lambda^{(n)} \quad (\text{if the limit exists})$$

(i.e. the long-term probability distribution between states).

As such, it has the following properties:

1. $\pi_i \geq 0 \quad \forall i \in S$ (probabilities are not negative)
2. $\sum_{i \in S} \pi_i = 1$ (total prob. of being anywhere)
3. $\begin{cases} \pi = \pi P \\ \pi P^n = \pi, \forall n. \end{cases}$ (1, or any number, further transition doesn't make a difference)

Theorem (Ergös-Feller-Pollard)

For all states i, j in an irreducible, aperiodic MC:

- 1.) if the chain is transient, $\lim_{n \rightarrow \infty} P_{ij}^{(n)} = 0$ (or $\lim_{n \rightarrow \infty} P^n = \underset{\text{zero-matrix}}{\sum}$)
- 2.) if the chain is recurrent, $\lim_{n \rightarrow \infty} P_{ij}^{(n)} = \pi_j$ (each row of P converges to π)
- 2.1.) if chain is null-recurrent, $\pi = 0$
- 2.2.) if chain is positive-recurrent, every $\pi_i > 0$, $\sum_{i \in S} \pi_i = 1$

In case (2), mean time of return to i :

$$M_i = \frac{1}{\pi_i}. \quad (\text{note: } M_i = \infty \text{ if } \pi_i = 0).$$

Proof: non-examinable.

Time-reversibility

We can define a reversed chain by:

$$Y_n = X_{-n} \quad \text{for } -\infty < n < \infty$$

(Y_n) is also an MC with same stationary distribution.

Definition:

A MC is reversible if the transition matrices of (X_n) and (Y_n) are equal.

Theorem

A MC is reversible iff $\pi_i P_{ij} = \pi_j P_{ji} \quad \forall i, j \in S$

Proof:

Since

$Q = (q_{ij})$ is the transition matrix of (Y_n) .

Then

$$\begin{aligned} q_{ij} &= P(Y_{n+1} = j | Y_n = i) = P(X_{-n-1} = j | X_{-n} = i) \\ &= \frac{P(X_{-n} = i | X_{-n-1} = j) \cdot P(X_{-n-1} = j)}{P(X_{-n} = i)} = \frac{P_{ji} \cdot \pi_j}{\pi_i} \quad (\text{by Bayes' theorem}) \end{aligned}$$

For MC to be reversible, we need

$$P_{ij} = \frac{P_{ji} \pi_j}{\pi_i} \quad \Rightarrow \quad P_{ij} \pi_i = P_{ji} \pi_j \quad \square$$

Theorem

For an irreducible MC, if there exists a vector π such that:

- 1) $0 \leq \pi_i \leq 1$ and $\sum_i \pi_i = 1$
- 2) $\pi_i P_{ij} = \pi_j P_{ji} \quad \forall i, j \in S$

then the MC is reversible

(local balance condition)

with stationary distribution π .

Proof:

$$(\pi P)_j = \sum_{i \in S} \pi_i P_{ij} = \sum_{i \in S} \pi_i P_{ji} = \pi_j \sum_{i \in S} P_{ji} = \pi_j$$

so $\pi P = \pi$ and distribution π is stationary. \square

Ergodic results

Let's consider the proportion of time we spend in a state i .
 If we visit state i $V_i(n)$ times before time n , then:

$$V_i(n) = \sum_{k=0}^n \mathbb{I}\{\bar{X}_k = i\}.$$

Then $\frac{V_i(n)}{n}$ is the proportion of time we spend in state i before ~~the~~ time n .

Theorem (Ergodic theorem)

For an irreducible MC (X_n) with transition matrix P

~~$$\lim_{n \rightarrow \infty} \frac{V_i(n)}{n} = \frac{1}{\mu_i} = \pi_i = 1.$$~~

where μ_i is expected return time to state i .

Proof: in the notes.

~~Markov chains~~

b.) Hidden Markov Models (HMMs)

Sometimes statistical model of observed data is constructed from an underlying but hidden MC. Can be useful in: OCR, Nat. language processing etc.

HMMs: We have an MC with transition matrix P , but the ~~state~~ current state of the MC is not directly observable. Instead we can observe a randomly chosen (distribution is given) token y_n . For each state there is a distribution b_i , such that $(b_i)_t = P(y_n=t | X_n=i)$.

Three central problems:

Evaluation

Given an observed sequence of tokens, determine its probability given the HMM parameters.

$$P(y_1=a_1, y_2=a_2, y_3=a_3, \dots | \text{HMM params})$$

In practice, solved by forward algorithm.

Decoding

Given a sequence of observed tokens and HMM parameters, determine the best-fitting sequence of hidden underlying states.

In practice, solved by Viterbi algorithm.

Learning

Given a sequence of tokens, determine the HMM parameters to maximize its probability.

$$P(y_1=a_1, \dots, y_n=a_n | \text{HMM params})$$

Solved by the Baum-Welch iterative method.