

# Mathematical Methods for Computer Science

## Probability methods

Dr R.J. Gibbens

Computer Laboratory  
University of Cambridge

Computer Science Tripos, Part IB  
Michaelmas Term 2014/15

Last revised: 2014-10-26 (980990d)

# Outline

- ▶ Probability methods (10 lectures, Dr R.J. Gibbens)
  - ▶ Probability generating functions (2 lectures)
  - ▶ Inequalities and limit theorems (3 lectures)
  - ▶ Stochastic processes (5 lectures)
- ▶ Fourier and related methods (6 lectures, Professor J. Daugman)

## Reference books (Probability methods)

- ▶ (\*) Ross, Sheldon M.  
*Probability Models for Computer Science.*  
Harcourt/Academic Press, 2002
- ▶ Mitzenmacher, Michael & Upfal, Eli.  
*Probability and Computing: Randomized Algorithms and Probabilistic Analysis.*  
Cambridge University Press, 2005

## Some notation

RV	random variable
IID	independent, identically distributed
PGF	probability generating function $G_X(z)$
MGF	moment generating function $M_X(t)$
$X \sim U(0, 1)$	RV $X$ has the distribution $U(0, 1)$ , etc
$\mathbb{I}(A)$	indicator function of the event $A$
$\mathbb{P}(A)$	probability that event $A$ occurs
$\mathbb{E}(X)$	expected value of RV $X$
$\mathbb{E}(X^n)$	$n^{th}$ moment of RV $X$ , for $n = 1, 2, \dots$
$F_X(x)$	distribution function, $F_X(x) = \mathbb{P}(X \leq x)$
$f_X(x)$	density of RV $X$ given, when it exists, by $F'_X(x)$

# Probability generating functions

# Probability generating functions (PGF)

A very common situation is when a RV,  $X$ , can take only non-negative integer values. For example,  $X$  may count the number of random events to occur in a fixed period of time. The probability mass function,  $\mathbb{P}(X = k)$ , is given by a sequence of values  $p_0, p_1, p_2, \dots$  where

$$p_k = \mathbb{P}(X = k) \geq 0 \quad \forall k \in \{0, 1, 2, \dots\} \quad \text{and} \quad \sum_{k=0}^{\infty} p_k = 1.$$

This sequence of terms can be “wrapped together” to define a function called the **probability generating function** (PGF) as follows.

## Definition (Probability generating function)

The **probability generating function**,  $G_X(z)$ , of a (non-negative integer-valued) RV  $X$  is defined as

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

for all values of  $z$  such that the sum converges.

# Elementary properties of the PGF

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

$$G_X(0) = p_0 \quad \text{and} \quad G_X(1) = 1.$$

2. If  $g(t) = z^t$  then

$$G_X(z) = \sum_{k=0}^{\infty} p_k z^k = \sum_{k=0}^{\infty} g(k) \mathbb{P}(X = k) = \mathbb{E}(g(X)) = \mathbb{E}(z^X).$$

3. The PGF is defined for all  $|z| \leq 1$  since

$$\sum_{k=0}^{\infty} |p_k z^k| \leq \sum_{k=0}^{\infty} p_k = 1.$$

4. Importantly, the PGF **characterizes** the distribution of a RV in the sense that

$$G_X(z) = G_Y(z) \quad \forall z$$

if and only if

$$\mathbb{P}(X = k) = \mathbb{P}(Y = k) \quad \forall k \in \{0, 1, 2, \dots\}.$$

# Examples of PGFs

## Example (Bernoulli distribution)

$$G_X(z) = q + pz \quad \text{where } q = 1 - p.$$

## Example (Binomial distribution, $\text{Bin}(n, p)$ )

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

## Example (Geometric distribution, $\text{Geo}(p)$ )

$$G_X(z) = \sum_{k=1}^{\infty} p q^{k-1} z^k = pz \sum_{k=0}^{\infty} (qz)^k = \frac{pz}{1 - qz} \text{ if } |z| < q^{-1} \text{ and } q = 1 - p.$$



# Examples of PGFs, ctd

**Example (Uniform distribution,  $U(1, n)$ )**

$$G_X(z) = \sum_{k=1}^n z^k \frac{1}{n} = \frac{z}{n} \sum_{k=0}^{n-1} z^k = \frac{z}{n} \frac{(1 - z^n)}{(1 - z)}.$$

**Example (Poisson distribution,  $\text{Pois}(\lambda)$ )**

$$G_X(z) = \sum_{k=0}^{\infty} \frac{\lambda^k e^{-\lambda}}{k!} z^k = e^{\lambda z} e^{-\lambda} = e^{\lambda(z-1)}.$$

## Derivatives of the PGF

We can derive a very useful property of the PGF by considering the derivative,  $G'_X(z)$ , with respect to  $z$ . Assume we can interchange the order of differentiation and summation, so that

$$\begin{aligned} G'_X(z) &= \frac{d}{dz} \left( \sum_{k=0}^{\infty} z^k \mathbb{P}(X = k) \right) \\ &= \sum_{k=0}^{\infty} \frac{d}{dz} (z^k) \mathbb{P}(X = k) \\ &= \sum_{k=0}^{\infty} k z^{k-1} \mathbb{P}(X = k) \end{aligned}$$

then putting  $z = 1$  we have that

$$G'_X(1) = \sum_{k=0}^{\infty} k \mathbb{P}(X = k) = \mathbb{E}(X)$$

the expectation of the RV  $X$ .

## Further derivatives of the PGF

Taking the second derivative gives

$$G_X''(z) = \sum_{k=0}^{\infty} k(k-1)z^{k-2}\mathbb{P}(X=k).$$

So that,

$$G_X''(1) = \sum_{k=0}^{\infty} k(k-1)\mathbb{P}(X=k) = \mathbb{E}(X(X-1))$$

Generally, we have the following result.

### Theorem

*If the RV  $X$  has PGF  $G_X(z)$  then the  $r^{\text{th}}$  derivative of the PGF, written  $G_X^{(r)}(z)$ , evaluated at  $z = 1$  is such that*

$$G_X^{(r)}(1) = \mathbb{E}(X(X-1)\cdots(X-r+1)).$$

## Using the PGF to calculate $\mathbb{E}(X)$ and $\text{Var}(X)$

We have that

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

and

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

For example, if  $X$  is a RV with the  $\text{Pois}(\lambda)$  distribution then  $G_X(z) = e^{\lambda(z-1)}$ . Thus,  $G'_X(z) = \lambda e^{\lambda(z-1)}$ ,  $G''_X(z) = \lambda^2 e^{\lambda(z-1)}$  and so  $G'_X(1) = \lambda$  and  $G''_X(1) = \lambda^2$ . So, finally,

$$\mathbb{E}(X) = \lambda \quad \text{and} \quad \text{Var}(X) = \lambda^2 + \lambda - \lambda^2 = \lambda.$$

# Sums of independent random variables

The following theorem shows how PGFs can be used to find the PGF of the sum of independent RVs.

## Theorem

If  $X$  and  $Y$  are *independent* RVs with PGFs  $G_X(z)$  and  $G_Y(z)$  respectively then

$$G_{X+Y}(z) = G_X(z)G_Y(z).$$

## Proof.

Using the independence of  $X$  and  $Y$  we have that

$$\begin{aligned} G_{X+Y}(z) &= \mathbb{E}(z^{X+Y}) \\ &= \mathbb{E}(z^X z^Y) \\ &= \mathbb{E}(z^X) \mathbb{E}(z^Y) \\ &= G_X(z) G_Y(z) \end{aligned}$$



## PGF example: Poisson RVs

For example, suppose that  $X$  and  $Y$  are independent RVs with  $X \sim \text{Pois}(\lambda_1)$  and  $Y \sim \text{Pois}(\lambda_2)$ , respectively.

Then

$$\begin{aligned} G_{X+Y}(z) &= G_X(z)G_Y(z) \\ &= e^{\lambda_1(z-1)} e^{\lambda_2(z-1)} \\ &= e^{(\lambda_1+\lambda_2)(z-1)}. \end{aligned}$$

Hence  $X + Y \sim \text{Pois}(\lambda_1 + \lambda_2)$  is again a Poisson RV but with the parameter  $\lambda_1 + \lambda_2$ .

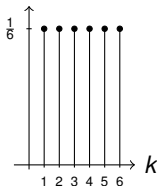
# PGF example: Uniform RVs



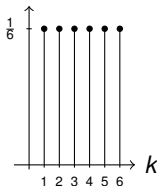
Consider the case of two fair dice with IID outcomes  $X$  and  $Y$ , respectively, so that  $X \sim U(1, 6)$  and  $Y \sim U(1, 6)$ . Let the total score be  $T = X + Y$  and consider the PGF of  $T$  given by  $G_T(z) = G_X(z)G_Y(z)$ . Then

$$\begin{aligned} G_T(z) &= \sum_{k=0}^{\infty} p_k z^k = \frac{1}{6}(z + z^2 + \cdots + z^6) \frac{1}{6}(z + z^2 + \cdots + z^6) \\ &= \frac{1}{36} [z^2 + 2z^3 + 3z^4 + 4z^5 + 5z^6 + 6z^7 + \\ &\quad 5z^8 + 4z^9 + 3z^{10} + 2z^{11} + z^{12}]. \end{aligned}$$

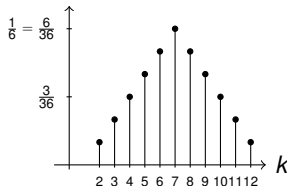
$\mathbb{P}(X = k)$



$\mathbb{P}(Y = k)$



$\mathbb{P}(T = k)$



# Limits and inequalities



# Limits and inequalities

We are familiar with limits of real numbers. For example, if  $x_n = 1/n$  for  $n = 1, 2, \dots$  then  $\lim_{n \rightarrow \infty} x_n = 0$  whereas if  $x_n = (-1)^n$  no such limit exists. Behaviour **in the long-run** or **on average** is an important characteristic of everyday life.

We will be concerned with these notions of limiting behaviour when the real numbers  $x_n$  are replaced by random variables  $X_n$ . As we shall see there are several distinct notions of convergence that can be considered.

To study these forms of convergence and the limiting theorems that emerge we shall also gather a very useful collection of concepts and tools for the probabilistic analysis of models, algorithms and systems.

# Probabilistic inequalities

To help assess how close RVs are to each other it is useful to have methods that provide upper bounds on probabilities of the form

$$\mathbb{P}(X \geq a)$$

for fixed constants  $a$ .

We shall consider several such bounds and related inequalities.

- ▶ Markov's inequality
- ▶ Chebyshev's inequality
- ▶ Chernoff's inequality

We will use  $\mathbb{I}(A)$  for the indicator RV which is 1 if  $A$  occurs and 0 otherwise. Observe that we have for such indicator RVs that

$$\mathbb{E}(\mathbb{I}(A)) = 0 \times \mathbb{P}(A^c) + 1 \times \mathbb{P}(A) = \mathbb{P}(A).$$

## Theorem (Markov's inequality)

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

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

### Proof.

We have that

$$\mathbb{I}(|X| \geq a) = \begin{cases} 1 & |X| \geq a \\ 0 & \text{otherwise.} \end{cases}$$

Clearly,

$$|X| \geq a \mathbb{I}(|X| \geq a)$$

hence

$$\mathbb{E}(|X|) \geq \mathbb{E}(a \mathbb{I}(|X| \geq a)) = a \mathbb{P}(|X| \geq a)$$

which yields the result.



## Theorem (Chebyshev's inequality)

Let  $X$  be a RV with mean  $\mu = \mathbb{E}(X)$  and finite variance  $\sigma^2 = \text{Var}(X)$  then for all  $a > 0$

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

### Proof.

Put  $Y = (X - \mu)^2 \geq 0$  then  $\mathbb{E}(Y) = \mathbb{E}((X - \mu)^2) = \text{Var}(X) = \sigma^2$ . So, by Markov's inequality, for all  $b > 0$

$$\mathbb{P}((X - \mu)^2 \geq b) = \mathbb{P}(Y \geq b) \leq \frac{\mathbb{E}(Y)}{b} = \frac{\sigma^2}{b}.$$

Now put  $b = a^2$  and noting that  $\mathbb{P}((X - \mu)^2 \geq a^2) = \mathbb{P}(|X - \mu| \geq a)$  we have that

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



# Moment generating function

## Definition

The **moment generating function** (MGF) of a RV  $X$ , written  $M_X(t)$ , is given by

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

and is defined for those values of  $t \in \mathbb{R}$  for which this expectation exists.

Using the power series  $e^x = 1 + x + x^2/2! + x^3/3! + \dots$  we see that

$$M_X(t) = \mathbb{E}(e^{tX}) = 1 + \mathbb{E}(X)t + \mathbb{E}(X^2)t^2/2! + \mathbb{E}(X^3)t^3/3! + \dots$$

and so the  $n^{\text{th}}$  moment of  $X$ ,  $\mathbb{E}(X^n)$ , is given by the coefficient of  $t^n/n!$  in the power series expansion of the MGF  $M_X(t)$ .

Note that for every RV,  $X$ , we have that  $M_X(0) = 1$  since

$$M_X(0) = \mathbb{E}(e^{0X}) = \mathbb{E}(1) = 1.$$

# Elementary properties of the MGF

1. If  $X$  has MGF  $M_X(t)$  then  $Y = aX + b$  has MGF  $M_Y(t) = e^{bt}M_X(at)$ .
2. If  $X$  and  $Y$  are **independent** then  $X + Y$  has MGF  $M_{X+Y}(t) = M_X(t)M_Y(t)$ .
3.  $\mathbb{E}(X^n) = M_X^{(n)}(0)$  where  $M_X^{(n)}$  is the  $n^{th}$  derivative of  $M_X$ .
4. If  $X$  is a discrete RV taking values  $0, 1, 2, \dots$  with PGF  $G_X(z) = \mathbb{E}(z^X)$  then  $M_X(t) = G_X(e^t)$ .

# Fundamental properties of the MGF

We will use without proof the following results.

1. **Uniqueness**: to each MGF there corresponds a unique distribution function having that MGF.  
In fact, if  $X$  and  $Y$  are RVs with the **same** MGF in some region  $-a < t < a$  where  $a > 0$  then  $X$  and  $Y$  have the **same** distribution.
2. **Continuity**: if distribution functions  $F_n(x)$  converge pointwise to a distribution function  $F(x)$ , the corresponding MGFs (where they exist) converge to the MGF of  $F(x)$ . Conversely, if a sequence of MGFs  $M_n(t)$  converge to  $M(t)$  which is continuous at  $t = 0$ , then  $M(t)$  is a MGF, and the corresponding distribution functions  $F_n(x)$  converge to the distribution function determined by  $M(t)$ .

## Example: exponential distribution

If  $X$  has an exponential distribution with parameter  $\lambda > 0$  then  $f_X(x) = \lambda e^{-\lambda x}$  for  $0 < x < \infty$ . Hence, for  $t < \lambda$ ,

$$\begin{aligned} M_X(t) &= \int_0^{\infty} e^{tx} \lambda e^{-\lambda x} dx = \int_0^{\infty} \lambda e^{-(\lambda-t)x} dx \\ &= \left[ -\frac{\lambda}{(\lambda-t)} e^{-(\lambda-t)x} \right]_0^{\infty} = \frac{\lambda}{\lambda-t}. \end{aligned}$$

For  $t < \lambda$

$$\frac{\lambda}{(\lambda-t)} = \left(1 - \frac{t}{\lambda}\right)^{-1} = 1 + \frac{t}{\lambda} + \frac{t^2}{\lambda^2} + \dots$$

and hence  $\mathbb{E}(X) = 1/\lambda$  and  $\mathbb{E}(X^2) = 2/\lambda^2$  so that

$$\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2 = 1/\lambda^2.$$



## Example: normal distribution

Consider a normal RV  $X \sim N(\mu, \sigma^2)$  then  $f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$  so that

$$\begin{aligned} M_X(t) &= \int_{-\infty}^{\infty} e^{tx} \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2} dx \\ &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(-2tx\sigma^2 + (x-\mu)^2)/2\sigma^2} dx. \end{aligned}$$

So, by completing the square,

$$\begin{aligned} M_X(t) &= e^{\mu t + \sigma^2 t^2/2} \left\{ \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(x-(\mu+t\sigma^2))^2/2\sigma^2} dx \right\} \\ &= e^{\mu t + \sigma^2 t^2/2}. \end{aligned}$$

## Example: uniform distribution

Consider a uniform RV  $X \sim U(a, b)$  for  $a < b$ . Then

$$f_X(x) = \begin{cases} \frac{1}{b-a} & a < x < b \\ 0 & \text{otherwise.} \end{cases}$$

Hence, for  $t \neq 0$ ,

$$\begin{aligned} M_X(t) &= \int_a^b \frac{e^{tx}}{b-a} dx \\ &= \left[ \frac{e^{tx}}{(b-a)t} \right]_a^b \\ &= \frac{e^{bt} - e^{at}}{(b-a)t}. \end{aligned}$$

and  $M_X(0) = 1$ .

## Theorem (Chernoff's bound)

Suppose that  $X$  has MGF  $M_X(t)$  and  $a \in \mathbb{R}$  then for all  $t > 0$

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

### Proof.

Using Markov's inequality, we have that since  $t > 0$

$$\begin{aligned}\mathbb{P}(X \geq a) &= \mathbb{P}(e^{tX} \geq e^{ta}) \\ &\leq \frac{\mathbb{E}(e^{tX})}{e^{ta}} \\ &= e^{-ta} M_X(t)\end{aligned}$$



Note that the above bound holds for all  $t > 0$  so we can select the **best** such bound by choosing  $t > 0$  to minimize  $e^{-ta} M_X(t)$ .

In fact, the upper bound also holds trivially if  $t = 0$  since the RHS is 1.

## Notions of convergence: $X_n \rightarrow X$ as $n \rightarrow \infty$

For a sequence of RVs  $(X_n)_{n \geq 1}$ , we shall define two distinct notions of convergence to some RV  $X$  as  $n \rightarrow \infty$ .

### Definition (Convergence in distribution)

$X_n \xrightarrow{D} X$  if  $F_{X_n}(x) \rightarrow F_X(x)$  for all points  $x$  at which  $F_X$  is continuous.

### Definition (Convergence in probability)

$X_n \xrightarrow{P} X$  if  $\mathbb{P}(|X_n - X| > \varepsilon) \rightarrow 0$  for all  $\varepsilon > 0$ .

There are further inter-related notions of convergence but two will suffice for our purposes.

## Theorem

If  $X_n \xrightarrow{P} X$  then  $X_n \xrightarrow{D} X$ .

## Proof

We prove this theorem as follows. Fix,  $\varepsilon > 0$  then

$$F_{X_n}(x) = \mathbb{P}(X_n \leq x \cap X > x + \varepsilon) + \mathbb{P}(X_n \leq x \cap X \leq x + \varepsilon)$$

since  $X > x + \varepsilon$  and  $X \leq x + \varepsilon$  form a partition. But if  $X_n \leq x$  and  $X > x + \varepsilon$  then  $|X_n - X| > \varepsilon$  and  $\{X_n \leq x \cap X \leq x + \varepsilon\} \subset \{X \leq x + \varepsilon\}$ . Therefore,

$$F_{X_n}(x) \leq \mathbb{P}(|X_n - X| > \varepsilon) + F_X(x + \varepsilon).$$

Similarly,

$$\begin{aligned} F_X(x - \varepsilon) &= \mathbb{P}(X \leq x - \varepsilon \cap X_n > x) + \mathbb{P}(X \leq x - \varepsilon \cap X_n \leq x) \\ &\leq \mathbb{P}(|X_n - X| > \varepsilon) + F_{X_n}(x). \end{aligned}$$

The proof is completed by noting that together these inequalities show that

$$F_X(x - \varepsilon) - \mathbb{P}(|X_n - X| > \varepsilon) \leq F_{X_n}(x) \leq \mathbb{P}(|X_n - X| > \varepsilon) + F_X(x + \varepsilon).$$

But  $X_n \xrightarrow{P} X$  implies that  $\mathbb{P}(|X_n - X| > \varepsilon) \rightarrow 0$ . So, as  $n \rightarrow \infty$ ,  $F_{X_n}(x)$  is squeezed between  $F_X(x - \varepsilon)$  and  $F_X(x + \varepsilon)$ .

Hence, if  $F_X$  is continuous at  $x$ ,  $F_{X_n}(x) \rightarrow F_X(x)$  and so  $X_n \xrightarrow{D} X$ . □

The converse does not hold in general. However, the problem sheet contains an exercise showing an important special case where the converse does hold.

# Limit theorems

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

$$S_n = X_1 + X_2 + \cdots + X_n \quad \text{and} \quad \bar{X}_n = S_n/n.$$

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

## Theorem (Weak Law of Large Numbers/WLLN)

*Suppose  $(X_n)_{n \geq 1}$  are IID RVs with finite mean  $\mu$  (and finite variance  $\sigma^2$ ) then  $\bar{X}_n \xrightarrow{P} \mu$ .*

Note that convergence to  $\mu$  in the WLLN (and SLLN) actually means convergence to a **degenerate** RV,  $X$ , with  $\mathbb{P}(X = \mu) = 1$ .

This is referred to as the weak law of large numbers since under more restrictive assumptions it holds (ie SLLN) for a stronger form of convergence known as **almost sure** convergence. Under the SLLN with almost sure convergence we would have that  $\mathbb{P}(\bar{X}_n \rightarrow \mu) = 1$ .

## Theorem (Weak Law of Large Numbers/WLLN)

Suppose  $(X_n)_{n \geq 1}$  are IID RVs with finite mean  $\mu$  and finite variance  $\sigma^2$  then  $\bar{X}_n \xrightarrow{P} \mu$ .

### Proof.

Recall that  $\mathbb{E}(\bar{X}_n) = \mu$  and  $\text{Var}(\bar{X}_n) = \sigma^2/n$ . Hence, by Chebyshev's inequality applied to  $\bar{X}_n$  for all  $\varepsilon > 0$

$$0 \leq \mathbb{P}(|\bar{X}_n - \mu| > \varepsilon) \leq \frac{\sigma^2/n}{\varepsilon^2} = \frac{\sigma^2}{n\varepsilon^2}$$

and so, letting  $n \rightarrow \infty$ ,

$$\mathbb{P}(|\bar{X}_n - \mu| > \varepsilon) \rightarrow 0$$

hence  $\bar{X}_n \xrightarrow{P} \mu$  as required. □



## Applications: estimating probabilities

Suppose we wish to estimate the probability,  $p$ , that we succeed when we play some game or perform some experiment. For  $i = 1, \dots, n$ , let

$$X_i = \mathbb{I}(\{i^{\text{th}} \text{ game is success}\}).$$

So  $\bar{X}_n = m/n$  if we succeed  $m$  times in  $n$  attempts.

We have that  $\mu = \mathbb{E}(X_i) = \mathbb{P}(X_i = 1) = p$  so then

$$m/n \xrightarrow{P} p$$

by the WLLN.

Thus we have shown the important result that the empirical estimate of the probability of some event by its observed sample frequency converges in probability to the correct but usually unknown value as the number of samples grows.

This result forms the basis of all simulation methods.

# Applications: Shannon's entropy

## Theorem (Asymptotic Equipartition Property/AEP)

If  $X_n$  is a sequence of IID discrete RV with probability distribution given by  $\mathbb{P}(X_i = x) = p(x)$  for each  $x \in I$  then

$$-\frac{1}{n} \log_2 p(X_1, X_2, \dots, X_n) \xrightarrow{P} H(X)$$

where Shannon's *entropy* is defined by

$$H(X) = H(X_1) = \dots = H(X_n) = - \sum_{x \in I} p(x) \log_2 p(x)$$

and

$$p(x_1, x_2, \dots, x_n) = \prod_{i=1}^n p(x_i)$$

is the joint probability distribution of the  $n$  IID RVs  $X_1, X_2, \dots, X_n$ .

### Proof.

Observe that  $p(X_i)$  is a RV taking the value  $p(x)$  with probability  $p(x)$  and similarly  $p(X_1, X_2, \dots, X_n)$  is a RV taking a value  $p(x_1, x_2, \dots, x_n)$  with probability  $p(x_1, x_2, \dots, x_n)$ . Therefore,

$$\begin{aligned} -\frac{1}{n} \log_2 p(X_1, X_2, \dots, X_n) &= -\frac{1}{n} \log_2 \prod_{i=1}^n p(X_i) \\ &= -\frac{1}{n} \sum_{i=1}^n \log_2 p(X_i) \\ &= \frac{1}{n} \sum_{i=1}^n (-\log_2 p(X_i)) \\ &\xrightarrow{P} \mathbb{E}(-\log_2 p(X_i)) && \text{by WLLN} \\ &= -\sum_{x \in I} p(x) \log_2 p(x) \\ &= H(X) \end{aligned}$$



# AEP implications

By the AEP, for all  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|-\frac{1}{n} \log_2 p(X_1, X_2, \dots, X_n) - H(X)| \leq \varepsilon) = 1$$

$$\lim_{n \rightarrow \infty} \mathbb{P}(H(X) - \varepsilon \leq -\frac{1}{n} \log_2 p(X_1, X_2, \dots, X_n) \leq H(X) + \varepsilon) = 1$$

$$\lim_{n \rightarrow \infty} \mathbb{P}(-n(H(X) - \varepsilon) \geq \log_2 p(X_1, X_2, \dots, X_n) \geq -n(H(X) + \varepsilon)) = 1$$

$$\lim_{n \rightarrow \infty} \mathbb{P}(2^{-n(H(X)+\varepsilon)} \leq p(X_1, X_2, \dots, X_n) \leq 2^{-n(H(X)-\varepsilon)}) = 1$$

Thus, the sequences of outcomes  $(x_1, x_2, \dots, x_n) \in A_n^\varepsilon$  where

$$A_n^\varepsilon = \{(x_1, x_2, \dots, x_n) : 2^{-n(H(X)+\varepsilon)} \leq p(x_1, x_2, \dots, x_n) \leq 2^{-n(H(X)-\varepsilon)}\}$$

have a high probability and are referred to as **typical sequences**. An efficient (optimal) coding is to assign short codewords to such sequences leaving longer codewords for any non-typical sequence. Such long codewords must arise only rarely in the limit and we would need around  $n(H(X) + \varepsilon)$  bits to distinguish these typical codewords.

# Central limit theorem

## Theorem (Central limit theorem/CLT)

*Let  $(X_n)_{n \geq 1}$  be a sequence of IID RVs with mean  $\mu$ , variance  $\sigma^2$  and whose moment generating function converges in some interval  $-a < t < a$  with  $a > 0$ . Then*

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

# Proof of CLT

Set  $Y_i = (X_i - \mu)/\sigma$  then  $\mathbb{E}(Y_i) = 0$  and  $\mathbb{E}(Y_i^2) = \text{Var}(Y_i) = 1$  so

$$M_{Y_i}(t) = 1 + \frac{t^2}{2} + o(t^2)$$

where  $o(t^2)$  refers to terms of higher order than  $t^2$  which will therefore tend to 0 as  $t \rightarrow 0$ . Also,

$$Z_n = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i.$$

Hence,

$$\begin{aligned} M_{Z_n}(t) &= \left( M_{Y_i} \left( \frac{t}{\sqrt{n}} \right) \right)^n \\ &= \left( 1 + \frac{t^2}{2n} + o \left( \frac{t^2}{n} \right) \right)^n \\ &\rightarrow e^{t^2/2} \quad \text{as} \quad n \rightarrow \infty. \end{aligned}$$

But  $e^{t^2/2}$  is the MGF of the  $N(0, 1)$  distribution so, together with the continuity property, the CLT now follows.

## CLT example

Suppose  $X_1, X_2, \dots, X_n$  are the IID RVs showing the  $n$  sample outcomes of a 6-sided die with common distribution

$$\mathbb{P}(X_i = j) = p_j, \quad j = 1, 2, \dots, 6$$

Set  $S_n = X_1 + X_2 + \dots + X_n$ , the total score obtained, and consider the two cases

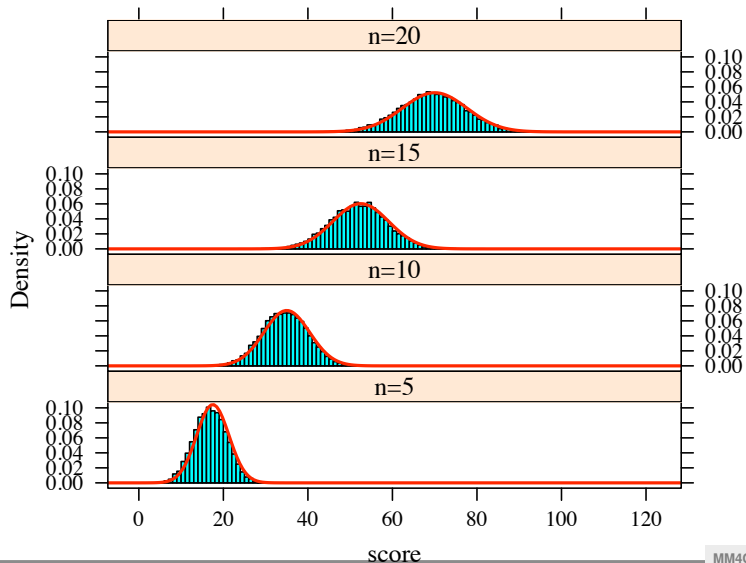
- ▶ **symmetric**:  $(p_j) = (1/6, 1/6, 1/6, 1/6, 1/6, 1/6)$  so that  $\mu = \mathbb{E}(X_i) = 3.5$  and  $\sigma^2 = \text{Var}(X_i) \approx 2.9$
- ▶ **asymmetric**:  $(p_j) = (0.2, 0.1, 0.0, 0.0, 0.3, 0.4)$  so that  $\mu = \mathbb{E}(X_i) = 4.3$  and  $\sigma^2 = \text{Var}(X_i) \approx 4.0$

for varying sample sizes  $n = 5, 10, 15$  and  $20$ .

The CLT tells us that for large  $n$ ,  $S_n$  is approximately distributed as  $N(n\mu, n\sigma^2)$  where  $\mu$  and  $\sigma^2$  are the mean and variance, respectively, of  $X_i$ .

# CLT example: symmetric

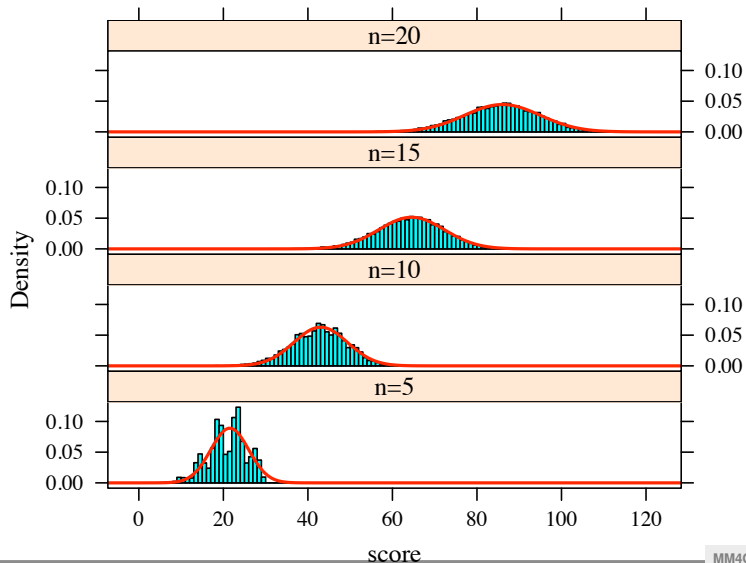
10,000 replications





# CLT example: asymmetric

10,000 replications



# Confidence intervals I

One of the major statistical applications of the CLT is to the construction of **confidence intervals**. The CLT shows that

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

is asymptotically distributed as  $N(0, 1)$ . If, the true value of  $\sigma^2$  is unknown we may estimate it by the **sample variance** given by

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

For instance, it can be shown that  $\mathbb{E}(S^2) = \sigma^2$  and then

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

is approximately distributed as  $N(0, 1)$  for large  $n$ .

## Confidence intervals II

Define  $z_\alpha$  so that  $\mathbb{P}(Z > z_\alpha) = \alpha$  where  $Z \sim N(0, 1)$  and so

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

Hence,

$$\begin{aligned}\mathbb{P}\left(-z_{\alpha/2} < \frac{\bar{X}_n - \mu}{S/\sqrt{n}} < z_{\alpha/2}\right) &\approx 1 - \alpha \\ \mathbb{P}\left(\bar{X}_n - z_{\alpha/2} \frac{S}{\sqrt{n}} < \mu < \bar{X}_n + z_{\alpha/2} \frac{S}{\sqrt{n}}\right) &\approx 1 - \alpha.\end{aligned}$$

The interval between the pair of end points  $\bar{X}_n \pm z_{\alpha/2} S/\sqrt{n}$  is thus an (approximate)  $100(1 - \alpha)$  percent **confidence interval** for the unknown parameter  $\mu$ .

## Confidence intervals: example

Consider a collection of  $n$  IID RVs,  $X_i$ , with common distribution  $X_i \sim \text{Pois}(\lambda)$ . Hence,

$$\mathbb{P}(X_i = j) = \frac{\lambda^j e^{-\lambda}}{j!} \quad j = 0, 1, \dots$$

with mean  $\mathbb{E}(X_i) = \lambda$ .

Then a 95% confidence interval for the (unknown) mean value  $\lambda$  is given by

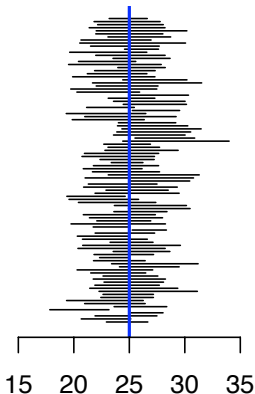
$$\bar{X}_n \pm 1.96S/\sqrt{n}$$

where  $z_{0.025} = 1.96$ .

Alternatively, to obtain 99% confidence intervals replace 1.96 by  $z_{0.005} = 2.58$  for a confidence interval  $\bar{X}_n \pm 2.58S/\sqrt{n}$ .

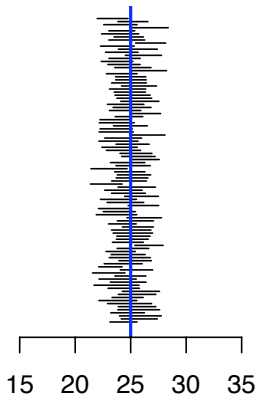
## 95% confidence intervals: illustration with $\lambda = 25$ and $\alpha = 5\%$

100 runs,  $n = 10$



confidence interval

100 runs,  $n = 40$



confidence interval

# Monte Carlo simulation and randomized algorithms

Suppose we wish to estimate the value of  $\pi$ . One way to proceed is to perform the following experiment. Select a point  $(X, Y) \in [-1, 1]^2$  with  $X$  and  $Y$  chosen independently and uniformly in  $[-1, 1]$ . Now consider those points within unit distance of the origin then

$$\mathbb{P}((X, Y) \text{ lies in unit circle}) = \mathbb{P}(X^2 + Y^2 \leq 1) = \frac{\text{area of circle}}{\text{area of square}} = \frac{\pi}{4}.$$

Suppose we have access to a stream of random variables  $U_i \sim U(0, 1)$  then  $2U_i - 1 \sim U(-1, 1)$ . Now set  $X_i = 2U_{2i-1} - 1$ ,  $Y_i = 2U_{2i} - 1$  and  $H_i = \mathbb{I}(\{X_i^2 + Y_i^2 \leq 1\})$  so that

$$\mathbb{E}(H_i) = \mathbb{P}(X_i^2 + Y_i^2 \leq 1) = \frac{\pi}{4}.$$

Hence by the WLLN the proportion of points  $(X_i, Y_i)$  falling within the unit circle converges in probability to  $\pi/4$ . Furthermore, the CLT can be used to form confidence intervals.

This a simple example of a randomized algorithm to solve a deterministic problem.

# Stochastic processes

# Random walks

Consider a sequence  $Y_1, Y_2, \dots$  of IID RVs with  $\mathbb{P}(Y_i = 1) = p$  and  $\mathbb{P}(Y_i = -1) = 1 - p$  with  $p \in [0, 1]$ .

## Definition (Simple random walk)

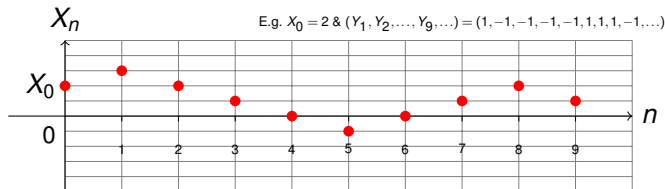
The **simple random walk** is a sequence of RVs  $\{X_n \mid n \in \{1, 2, \dots\}\}$  defined by

$$X_n = X_0 + Y_1 + Y_2 + \dots + Y_n$$

where  $X_0 \in \mathbb{R}$  is the starting value.

## Definition (Simple symmetric random walk)

A **simple symmetric random walk** is a simple random walk with the choice  $p = 1/2$ .





# Examples



Practical examples of random walks abound across the physical sciences (motion of atomic particles) and the non-physical sciences (epidemics, gambling, asset prices, cryptocurrencies).

The following is a simple model for the operation of a casino. Suppose that a gambler enters with a capital of  $\pounds X_0$ . At each stage the gambler places a stake of  $\pounds 1$  and with probability  $p$  wins the gamble otherwise the stake is lost. If the gambler wins the stake is returned together with an additional sum of  $\pounds 1$ .

Thus at each stage the gambler's capital increases by  $\pounds 1$  with probability  $p$  or decreases by  $\pounds 1$  with probability  $1 - p$ .

The gambler's capital  $X_n$  at stage  $n$  thus follows a simple random walk **except** that the gambler is **bankrupt** if  $X_n$  reaches  $\pounds 0$  and then can not continue to any further stages.

# Returning to the starting state for a simple random walk

Let  $X_n$  be a simple random walk and

$$r_n = \mathbb{P}(X_n = X_0) \quad \text{for } n = 1, 2, \dots$$

the probability of returning to the starting state at time  $n$ .  
We will show the following theorem.

## Theorem

*If  $n$  is odd then  $r_n = 0$  else if  $n = 2m$  is even then*

$$r_n = \binom{2m}{m} p^m (1-p)^m.$$

### Proof.

The position of the random walk will change by an amount

$$X_n - X_0 = Y_1 + Y_2 + \cdots + Y_n$$

between times 0 and  $n$ . Hence, for this change  $X_n - X_0$  to be 0 there must be an equal number of up steps as down steps. This can never happen if  $n$  is odd and so  $r_n = 0$  in this case. If  $n = 2m$  is even then note that the number of up steps in a total of  $n$  steps is a binomial RV with parameters  $2m$  and  $p$ . Thus,

$$r_n = \mathbb{P}(X_n - X_0 = 0) = \binom{2m}{m} p^m (1-p)^m.$$



This result tells us about the probability of returning to the starting state at a given time  $n$ .

We will now look at the probability that we ever return to our starting state. For convenience, and without loss of generality, we shall take our starting value as  $X_0 = 0$  from now on.

# Recurrence and transience of simple random walks

Note first that  $\mathbb{E}(Y_i) = p - (1 - p) = 2p - 1$  for each  $i \in \{1, 2, \dots\}$ . Thus there is a net drift upwards if  $p > 1/2$  and a net drift downwards if  $p < 1/2$ . Only in the case  $p = 1/2$  is there no net drift.

We say that the simple random walk is **recurrent** if it is certain to revisit its starting state at some time in the future and **transient** otherwise.

We shall prove the following theorem.

## Theorem

*For a simple random walk with starting state  $X_0 = 0$  the probability of revisiting the starting state is*

$$\mathbb{P}(X_n = 0 \text{ for some } n \in \{1, 2, \dots\}) = 1 - |2p - 1|.$$

Thus a simple random walk is **recurrent** only when  $p = 1/2$  and it is **transient** for all  $p \neq 1/2$ .

## Proof

Since we have assumed that  $X_0 = 0$  the event  $R_n = \{X_n = 0\}$  indicates that the simple random walk returns to its starting state at time  $n$ . Consider the event

$$F_n = \{X_n = 0, X_m \neq 0 \text{ for } m \in \{1, 2, \dots, (n-1)\}\}$$

that the random walk first revisits its starting state at time  $n$ . If  $R_n$  occurs then exactly one of  $F_1, F_2, \dots, F_n$  occurs. So,

$$\mathbb{P}(R_n) = \sum_{m=1}^n \mathbb{P}(R_n \cap F_m)$$

but

$$\mathbb{P}(R_n \cap F_m) = \mathbb{P}(F_m) \mathbb{P}(R_{n-m}) \quad \text{for } m \in \{1, 2, \dots, n\}$$

since we must first return at time  $m$  and then return a time  $n - m$  later which are independent events. So if we write  $f_n = \mathbb{P}(F_n)$  and  $r_n = \mathbb{P}(R_n)$  then

$$r_n = \sum_{m=1}^n f_m r_{n-m}.$$

Given the expression for  $r_n$  we now wish to solve these equations for  $f_m$ .

## Proof, ctd

Define generating functions for the sequences  $r_n$  and  $f_n$  by

$$R(z) = \sum_{n=0}^{\infty} r_n z^n \quad \text{and} \quad F(z) = \sum_{n=0}^{\infty} f_n z^n$$

where  $r_0 = 1$  and  $f_0 = 0$  and take  $|z| < 1$ . We have that

$$\begin{aligned} \sum_{n=1}^{\infty} r_n z^n &= \sum_{n=1}^{\infty} \sum_{m=1}^n f_m r_{n-m} z^n \\ &= \sum_{m=1}^{\infty} \sum_{n=m}^{\infty} f_m z^m r_{n-m} z^{n-m} \\ &= \sum_{m=1}^{\infty} f_m z^m \sum_{k=0}^{\infty} r_k z^k \\ &= F(z)R(z). \end{aligned}$$

The left hand side is  $R(z) - r_0 z^0 = R(z) - 1$  thus we have that

$$R(z) = R(z)F(z) + 1 \quad \text{if } |z| < 1.$$

## Proof, ctd

Now,

$$\begin{aligned} R(z) &= \sum_{n=0}^{\infty} r_n z^n \\ &= \sum_{m=0}^{\infty} r_{2m} z^{2m} \quad \text{as } r_n = 0 \text{ if } n \text{ is odd} \\ &= \sum_{m=0}^{\infty} \binom{2m}{m} (p(1-p)z^2)^m \\ &= (1 - 4p(1-p)z^2)^{-\frac{1}{2}}. \end{aligned}$$

The last step follows from the binomial series expansion of  $(1 - 4\theta)^{-\frac{1}{2}}$  and the choice  $\theta = p(1-p)z^2$ .

Hence,

$$F(z) = 1 - (1 - 4p(1-p)z^2)^{\frac{1}{2}} \quad \text{for } |z| < 1.$$

## Proof, ctd

But now

$$\begin{aligned}\mathbb{P}(X_n = 0 \text{ for some } n = 1, 2, \dots) &= \mathbb{P}(F_1 \cup F_2 \cup \dots) \\ &= f_1 + f_2 + \dots \\ &= \lim_{z \uparrow 1} \sum_{n=1}^{\infty} f_n z^n \\ &= F(1) \\ &= 1 - (1 - 4p(1-p))^{\frac{1}{2}} \\ &= 1 - ((2p-1)^2)^{\frac{1}{2}} \\ &= 1 - |2p-1|.\end{aligned}$$

So, finally, the simple random walk is certain to revisit its starting state just when  $p = 1/2$ .



## Mean return time

Consider the recurrent case when  $p = 1/2$  and set

$$T = \min\{n \geq 1 \mid X_n = 0\} \quad \text{so that} \quad \mathbb{P}(T = n) = f_n$$

where  $T$  is the time of the first return to the starting state. Then

$$\begin{aligned}\mathbb{E}(T) &= \sum_{n=1}^{\infty} n f_n \\ &= G'_T(1)\end{aligned}$$

where  $G_T(z)$  is the PGF of the RV  $T$  and for  $p = 1/2$  we have that  $4p(1-p) = 1$  so

$$G_T(z) = 1 - (1 - z^2)^{\frac{1}{2}}$$

so that

$$G'_T(z) = z(1 - z^2)^{-\frac{1}{2}} \rightarrow \infty \quad \text{as } z \uparrow 1.$$

Thus, the simple symmetric random walk ( $p = 1/2$ ) is recurrent but the expected time to first return to the starting state is **infinite**.

# The Gambler's ruin problem

We now consider a variant of the simple random walk. Consider two players A and B with a joint capital between them of  $\pounds N$ . Suppose that initially A has  $X_0 = \pounds a$  ( $0 \leq a \leq N$ ).

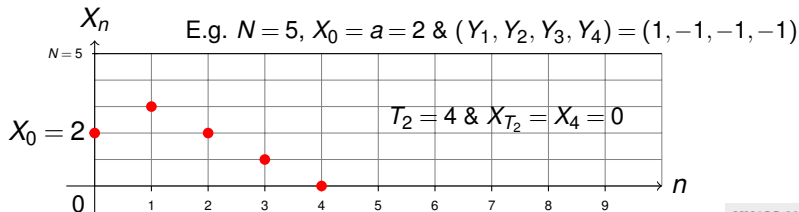
At each time step player B gives A  $\pounds 1$  with probability  $p$  and with probability  $q = (1 - p)$  player A gives  $\pounds 1$  to B instead. The outcomes at each time step are independent.

The game ends at the first time  $T_a$  if either  $X_{T_a} = \pounds 0$  or  $X_{T_a} = \pounds N$  for some  $T_a \in \{0, 1, \dots\}$ .

We can think of A's wealth,  $X_n$ , at time  $n$  as a simple random walk on the states  $\{0, 1, \dots, N\}$  with absorbing barriers at 0 and  $N$ .

Define the probability of ruin,  $\rho_a$ , for gambler A as

$$\rho_a = \mathbb{P}(\text{A is ruined}) = \mathbb{P}(\text{B wins}) \quad \text{for } 0 \leq a \leq N.$$



# Solution of the Gambler's ruin problem

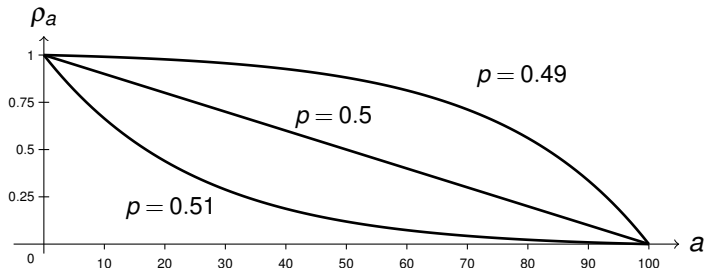
## Theorem

The probability of ruin when  $A$  starts with an initial capital of  $a$  is given by

$$\rho_a = \begin{cases} \frac{\theta^a - \theta^N}{1 - \theta^N} & \text{if } p \neq q \\ 1 - \frac{a}{N} & \text{if } p = q = 1/2 \end{cases}$$

where  $\theta = q/p$ .

For illustration here is a set of graphs of  $\rho_a$  for  $N = 100$  and three possible choices of  $p$ .



## Proof

Consider what happens at the first time step

$$\begin{aligned}\rho_a &= \mathbb{P}(\text{ruin} \cap Y_1 = +1 | X_0 = a) + \mathbb{P}(\text{ruin} \cap Y_1 = -1 | X_0 = a) \\ &= p\mathbb{P}(\text{ruin} | X_0 = a+1) + q\mathbb{P}(\text{ruin} | X_0 = a-1) \\ &= p\rho_{a+1} + q\rho_{a-1}\end{aligned}$$

Now look for a solution to this difference equation of the form  $\lambda^a$  with boundary conditions  $\rho_0 = 1$  and  $\rho_N = 0$ .

Try a solution of the form  $\rho_a = \lambda^a$  to give

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

Hence,

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

with solutions  $\lambda = 1$  and  $\lambda = q/p$ .

## Proof, ctd

If  $p \neq q$  there are two distinct solutions and the general solution of the difference equation is of the form  $A + B(q/p)^a$ .

Applying the boundary conditions

$$1 = \rho_0 = A + B \quad \text{and} \quad 0 = \rho_N = A + B(q/p)^N$$

we get

$$A = -B(q/p)^N$$

and

$$1 = B - B(q/p)^N$$

so

$$B = \frac{1}{1 - (q/p)^N} \quad \text{and} \quad A = \frac{-(q/p)^N}{1 - (q/p)^N}.$$

Hence,

$$\rho_a = \frac{(q/p)^a - (q/p)^N}{1 - (q/p)^N}.$$

## Proof, ctd

If  $p = q = 1/2$  then the general solution is  $C + Da$ .

So with the boundary conditions

$$1 = \rho_0 = C + D(0) \quad \text{and} \quad 0 = \rho_N = C + D(N).$$

Therefore,

$$C = 1 \quad \text{and} \quad 0 = 1 + D(N)$$

so

$$D = -1/N$$

and

$$\rho_a = 1 - a/N.$$

# Mean duration time

Set  $T_a$  as the time to be absorbed at either 0 or  $N$  starting from the initial state  $a$  and write  $\mu_a = \mathbb{E}(T_a)$ .

Then, conditioning on the first step as before

$$\mu_a = 1 + p\mu_{a+1} + q\mu_{a-1} \quad \text{for } 1 \leq a \leq N-1$$

and  $\mu_0 = \mu_N = 0$ .

It can be shown that  $\mu_a$  is given by

$$\mu_a = \begin{cases} \frac{1}{p-q} \left( N \frac{(q/p)^a - 1}{(q/p)^N - 1} - a \right) & \text{if } p \neq q \\ a(N-a) & \text{if } p = q = 1/2. \end{cases}$$

We skip the proof here but note the following cases can be used to establish the result.

**Case  $p \neq q$ :** trying a particular solution of the form  $\mu_a = ca$  shows that  $c = 1/(q-p)$  and the general solution is then of the form  $\mu_a = A + B(q/p)^a + a/(q-p)$ . Fixing the boundary conditions gives the result.

**Case  $p = q = 1/2$ :** now the particular solution is  $-a^2$  so the general solution is of the form  $\mu_a = A + Ba - a^2$  and fixing the boundary conditions gives the result.

# Markov chains

## Definition (Markov chain)

Suppose that  $(X_n)_{n \geq 0}$  is a sequence of discrete random variables taking values in some countable state space  $S$ . The sequence  $(X_n)$  is a **Markov chain** (MC) if

$$\mathbb{P}(X_n = x_n | X_0 = x_0, X_1 = x_1, \dots, X_{n-1} = x_{n-1}) = \mathbb{P}(X_n = x_n | X_{n-1} = x_{n-1})$$

for all  $n \geq 1$  and for all  $x_0, x_1, \dots, x_n \in S$ .

Since,  $S$  is countable we can always choose to label the possible values of  $X_n$  by integers and say that when  $X_n = i$  the Markov chain is in the “ **$i^{\text{th}}$  state at the  $n^{\text{th}}$  step**” or “**visits  $i$  at time  $n$** ”.



# Transition probabilities

The dynamics of the Markov chain are governed by the **transition probabilities**  $\mathbb{P}(X_n = j | X_{n-1} = i)$ .

## Definition (time-homogeneous MC)

A Markov chain  $(X_n)$  is **time-homogeneous** if

$$\mathbb{P}(X_n = j | X_{n-1} = i) = \mathbb{P}(X_1 = j | X_0 = i)$$

for all  $n \geq 1$  and states  $i, j \in S$ .

- ▶ We shall assume that our MCs are time-homogeneous unless explicitly stated otherwise.

# Transition matrix

## Definition (Transition matrix)

The **transition matrix**,  $P$ , of a MC  $(X_n)$  is given by  $P = (p_{ij})$  where for all  $i, j \in S$

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

- ▶ Note that  $P$  is a **stochastic matrix**, that is, it has non-negative entries ( $p_{ij} \geq 0$ ) and the row sums all equal one ( $\sum_j p_{ij} = 1$ ).
- ▶ The transition matrix completely characterizes the dynamics of the MC.

## Example

Suppose the states of the MC are  $S = \{1, 2, 3\}$  and that the transition matrix is given by

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

- ▶ Thus, in state 1 we are equally likely to be in any of the three states at the next step.
- ▶ In state 2, we can move with equal probabilities to 1 or 3 at the next step.
- ▶ Finally in state 3, we either move to state 1 with probability  $2/3$  or remain in state 3 at the next step.

# $n$ -step transition matrix

## Definition ( $n$ -step transition matrix)

The  $n$ -step transition matrix is  $P^{(n)} = (p_{ij}^{(n)})$  where

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

Thus  $P^{(1)} = P$  and we also set  $P^{(0)} = I$ , the  $|S| \times |S|$ -identity matrix.

# Chapman-Kolmogorov equations

## Theorem (Chapman-Kolmogorov)

For all states  $i, j$  and for all steps  $m, n$

$$p_{ij}^{(m+n)} = \sum_k p_{ik}^{(m)} p_{kj}^{(n)}.$$

Hence,  $P^{(m+n)} = P^{(m)} P^{(n)}$  and  $P^{(n)} = P^n$ , the  $n^{\text{th}}$  power of  $P$ .

**Proof.**

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

The Chapman-Kolmogorov equations tell us how the long-term evolution of the MC depends on the short-term evolution specified by the transition matrix.

If we let  $\lambda_i^{(n)} = \mathbb{P}(X_n = i)$  be the elements of a row vector  $\lambda^{(n)}$  specifying the distribution of the MC at the  $n^{\text{th}}$  time step then the follow holds.

### Lemma

*If  $m, n$  are non-negative integers then  $\lambda^{(m+n)} = \lambda^{(m)} P^{(n)}$  and so, in particular, if  $m = 0$*

$$\lambda^{(n)} = \lambda^{(0)} P^{(n)}$$

*where  $\lambda^{(0)}$  is the initial distribution  $\lambda_i^{(0)} = \mathbb{P}(X_0 = i)$  of the MC.*

### Proof.

$$\begin{aligned}\lambda_j^{(m+n)} &= \mathbb{P}(X_{m+n} = j) = \sum_i \mathbb{P}(X_{m+n} = j | X_m = i) \mathbb{P}(X_m = i) \\ &= \sum_i \lambda_i^{(m)} p_{ij}^{(n)} = \left( \lambda^{(m)} P^{(n)} \right)_j\end{aligned}$$

# Classification of states

## Definition (Accessibility)

If, for some  $n \geq 0$ ,  $p_{ij}^{(n)} > 0$  then we say that state  $j$  is **accessible** from state  $i$ , written  $i \rightsquigarrow j$ .

If  $i \rightsquigarrow j$  and  $j \rightsquigarrow i$  then we say that  $i$  and  $j$  **communicate**, written  $i \longleftrightarrow j$ .

Observe that the relation **communicates**  $\longleftrightarrow$  is

- ▶ reflexive
- ▶ symmetric
- ▶ transitive

and hence is an equivalence relation. The corresponding equivalence classes partition the state space into subsets of states, called **communicating classes**.

# Irreducibility

- ▶ A communicating class,  $C$ , that once entered can not be left is called **closed**, that is  $p_{ij} = 0$  for all  $i \in C, j \notin C$ .
- ▶ A closed communicating class consisting of a single state is called **absorbing**.
- ▶ When the state space forms a single communicating class, the MC is called **irreducible** and is called **reducible** otherwise.



# Recurrence and transience of MCs

Write for  $n \geq 1$

$$f_{ij}^{(n)} = \mathbb{P}(X_1 \neq j, \dots, X_{n-1} \neq j, X_n = j | X_0 = i)$$

so that  $f_{ij}^{(n)}$  is the probability starting in state  $i$  that we visit state  $j$  for the **first time** at time  $n$ . Also, let

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

the probability that we ever visit state  $j$ , starting in state  $i$ .

## Definition

- ▶ If  $f_{ii} < 1$  then state  $i$  is **transient**
- ▶ If  $f_{ii} = 1$  then state  $i$  is **recurrent**.

## Recurrence and transience, ctd

- ▶ Observe that if we return to a state  $i$  at some time  $n$  then the evolution of the MC is independent of the path before time  $n$ . Hence, the probability that we will return at least  $N$  times is  $f_{ii}^N$ .
- ▶ Now, if  $i$  is recurrent  $f_{ii}^N = 1$  for all  $N$  and we are sure to return to state  $i$  infinitely often.
- ▶ Conversely, if state  $i$  is transient then  $f_{ii}^N \rightarrow 0$  as  $N \rightarrow \infty$  and so there is zero probability of returning infinitely often.

## Theorem

- ▶  $i$  is transient  $\Leftrightarrow \sum_{n \geq 1} p_{ii}^{(n)}$  converges
- ▶  $i$  is recurrent  $\Leftrightarrow \sum_{n \geq 1} p_{ii}^{(n)}$  diverges

If  $i$  and  $j$  belong to the same communicating class then they are either both recurrent or both transient — the *solidarity property*.

## Proof

First, define generating functions

$$P_{ii}(z) = \sum_{n=0}^{\infty} p_{ii}^{(n)} z^n \quad \text{and} \quad F_{ii}(z) = \sum_{n=0}^{\infty} f_{ii}^{(n)} z^n$$

where we take  $p_{ii}^{(0)} = 1$  and  $f_{ii}^{(0)} = 0$ .

By examining the first time,  $r$ , that we return to  $i$ , we have for  $m = 1, 2, \dots$  that

$$p_{ii}^{(m)} = \sum_{r=1}^m f_{ii}^{(r)} p_{ii}^{(m-r)}.$$

Now multiply by  $z^m$  and summing over  $m$  we get

$$\begin{aligned} P_{ii}(z) &= 1 + \sum_{m=1}^{\infty} z^m p_{ii}^{(m)} \\ &= 1 + \sum_{m=1}^{\infty} z^m \sum_{r=1}^m f_{ii}^{(r)} p_{ii}^{(m-r)} \\ &= 1 + \sum_{r=1}^{\infty} f_{ii}^{(r)} z^r \sum_{m=r}^{\infty} p_{ii}^{(m-r)} z^{m-r} \\ &= 1 + F_{ii}(z) P_{ii}(z) \end{aligned}$$

Thus,  $P_{ii}(z) = 1/(1 - F_{ii}(z))$ . Now let  $z \nearrow 1$  then  $F_{ii}(z) \rightarrow F_{ii}(1) = f_{ii}$  and  $P_{ii}(z) \rightarrow \sum_n p_{ii}^{(n)}$ .

If  $i$  is transient then  $f_{ii} < 1$  so  $\sum_n p_{ii}^{(n)}$  converges. Conversely, if  $i$  is recurrent then  $f_{ii} = 1$  and  $\sum_n p_{ii}^{(n)}$  diverges.

Furthermore, if  $i$  and  $j$  are in the same class then there exist  $m$  and  $n$  so that  $p_{ij}^{(m)} > 0$  and  $p_{ji}^{(n)} > 0$ . Now, for all  $r \geq 0$

$$p_{ii}^{(m+r+n)} \geq p_{ij}^{(m)} p_{jj}^{(r)} p_{ji}^{(n)}$$

so that  $\sum_r p_{jj}^{(r)}$  and  $\sum_k p_{ii}^{(k)}$  diverge or converge together. □

# Mean recurrence time

First, let

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

be the time of the first visit to state  $j$  and set  $T_j = \infty$  if no such visit ever occurs.

Thus,  $\mathbb{P}(T_j = \infty | X_0 = i) > 0$  if and only if  $i$  is transient in which case  $\mathbb{E}(T_j | X_0 = i) = \infty$ .

## Definition (Mean recurrence time)

The **mean recurrent time**,  $\mu_i$ , of a state  $i$  is defined as

$$\mu_i = \mathbb{E}(T_i | X_0 = i) = \begin{cases} \sum_n n f_{ii}^{(n)} & \text{if } i \text{ is recurrent} \\ \infty & \text{if } i \text{ is transient.} \end{cases}$$

- Note that  $\mu_i$  may still be infinite when  $i$  is recurrent.

# Positive and null recurrence

## Definition

A recurrent state  $i$  is

- ▶ **positive recurrent** if  $\mu_i < \infty$  and
- ▶ **null recurrent** if  $\mu_i = \infty$ .

## Example: simple random walk

Recall the **simple random walk** where  $X_n = \sum_{i=1}^n Y_i$  where  $(Y_n)$  are IID RVs with  $\mathbb{P}(Y_i = 1) = p = 1 - \mathbb{P}(Y_i = -1)$ . Thus  $X_n$  is the position after  $n$  steps where we take unit steps up or down with probabilities  $p$  and  $1 - p$ , respectively.

It is clear that return to the origin is only possible after an even number of steps. Thus the sequence  $(p_{00}^{(n)})$  alternates between zero and a positive value.



# Periodicity

Let  $d_i$  be the greatest common divisor of  $\{n : p_{ii}^{(n)} > 0\}$ .

## Definition

- ▶ If  $d_i = 1$  then  $i$  is **aperiodic**.
- ▶ If  $d_i > 1$  then  $i$  is **periodic** with period  $d_i$ .
- ▶ It may be shown that the period is a class property, that is, if  $i, j \in C$  then  $d_i = d_j$ .

We will now concentrate on irreducible and aperiodic Markov chains.

# Stationary distributions

## Definition

The vector  $\pi = (\pi_j; j \in S)$  is a **stationary distribution** for the MC with transition matrix  $P$  if

1.  $\pi_j \geq 0$  for all  $j \in S$  and  $\sum_{j \in S} \pi_j = 1$
2.  $\pi = \pi P$ , or equivalently,  $\pi_j = \sum_{i \in S} \pi_i p_{ij}$ .

Such a distribution is stationary in the sense that  $\pi P^2 = (\pi P)P = \pi P = \pi$  and for all  $n \geq 0$

$$\pi P^n = \pi.$$

Thus if  $X_0$  has distribution  $\pi$  then  $X_n$  has distribution  $\pi$  for all  $n$ . Moreover,  $\pi$  is the **limiting distribution** of  $X_n$  as  $n \rightarrow \infty$ .

## Markov's example

Markov was lead to the notion of a Markov chain by study the patterns of vowels and consonants in text. In his original example, he found a transition matrix for the states {vowel, consonant) as

$$P = \begin{pmatrix} 0.128 & 0.872 \\ 0.663 & 0.337 \end{pmatrix}.$$

Taking successive powers of  $P$  we find

$$P^2 = \begin{pmatrix} 0.595 & 0.405 \\ 0.308 & 0.692 \end{pmatrix} \quad P^3 = \begin{pmatrix} 0.345 & 0.655 \\ 0.498 & 0.502 \end{pmatrix} \quad P^4 = \begin{pmatrix} 0.478 & 0.522 \\ 0.397 & 0.603 \end{pmatrix}.$$

As  $n \rightarrow \infty$ ,

$$P^n \rightarrow \begin{pmatrix} 0.432 & 0.568 \\ 0.432 & 0.568 \end{pmatrix}.$$

Check that  $\pi = (0.432, 0.568)$  is a stationary distribution, that is  $\pi P = \pi$ .

# Limiting behaviour as $n \rightarrow \infty$

## Theorem (Erdős-Feller-Pollard)

For all states  $i$  and  $j$  in an irreducible, aperiodic MC,

1. if the chain is transient,  $p_{ij}^{(n)} \rightarrow 0$
2. if the chain is recurrent,  $p_{ij}^{(n)} \rightarrow \pi_j$ , where
  - 2.1 (null recurrent) either, every  $\pi_j = 0$
  - 2.2 (positive recurrent) or, every  $\pi_j > 0$ ,  $\sum_j \pi_j = 1$  and  $\pi$  is the unique probability distribution solving  $\pi P = \pi$ .
3. In case (2), let  $T_i$  be the time to return to  $i$  then  $\mu_i = \mathbb{E}(T_i) = 1/\pi_i$  with  $\mu_i = \infty$  if  $\pi_i = 0$ .

## Proof.

Omitted.



## Remarks

- ▶ The limiting distribution,  $\pi$ , is seen to be a stationary one. Suppose the current distribution is given by  $\pi$  and consider the evolution of the MC for a further period of  $T$  steps. Since  $\pi$  is stationary, the probability of being in any state  $i$  remains  $\pi_i$ , so we will make around  $T\pi_i$  visits to  $i$ . Consequently, the mean time between visits to  $i$  would be  $T/(T\pi_i) = 1/\pi_i$ .
- ▶ Using  $\lambda_j^{(n)} = \mathbb{P}(X_n = j)$  and since  $\lambda^{(n)} = \lambda^{(0)}P^n$ 
  1. for transient or null recurrent states  $\lambda^{(n)} \rightarrow 0$ , that is,  $\mathbb{P}(X_n = j) \rightarrow 0$  for all states  $j$
  2. for a positive recurrent state,  $p^{(n)} \rightarrow \pi > 0$ , that is,  $\mathbb{P}(X_n = j) \rightarrow \pi_j > 0$  for all  $j$ , where  $\pi$  is the unique probability vector solving  $\pi P = \pi$ .
- ▶ Note the distinction between a transient and a null recurrent chain is that in a transient chain we might never make a return visit to some state  $i$  and there is zero probability that we will return infinitely often. However, in a null recurrent chain we are sure to make infinitely many return visits but the mean time between consecutive visits is infinite.

# Time-reversibility

Suppose now that  $(X_n : -\infty < n < \infty)$  is an irreducible, positive recurrent MC with transition matrix  $P$  and unique stationary distribution  $\pi$ . Suppose also that  $X_n$  has the distribution  $\pi$  for all  $-\infty < n < \infty$ . Now define the **reversed chain** by

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

Then  $(Y_n)$  is also a MC and where  $Y_n$  has the distribution  $\pi$ .

## Definition (Reversibility)

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

## Theorem

A MC  $(X_n)$  is reversible if and only if

$$\pi_i p_{ij} = \pi_j p_{ji} \quad \text{for all } i, j \in S.$$

## Proof.

Consider the transition probabilities  $q_{ij}$  of the MC  $(Y_n)$  then

$$\begin{aligned} q_{ij} &= \mathbb{P}(Y_{n+1} = j | Y_n = i) \\ &= \mathbb{P}(X_{-n-1} = j | X_{-n} = i) \\ &= \mathbb{P}(X_m = i | X_{m-1} = j) \mathbb{P}(X_{m-1} = j) / \mathbb{P}(X_m = i) \quad \text{where } m = -n \\ &= p_{ji} \pi_j / \pi_i. \end{aligned}$$

Hence,  $p_{ij} = q_{ij}$  if and only if  $\pi_i p_{ij} = \pi_j p_{ji}$ . □

## Theorem

For an irreducible chain, if there exists a vector  $\pi$  such that

1.  $0 \leq \pi_i \leq 1$  and  $\sum_i \pi_i = 1$
2.  $\pi_i p_{ij} = \pi_j p_{ji}$  for all  $i, j \in S$

then the MC is reversible with stationary distribution  $\pi$ .

## Proof.

Suppose that  $\pi$  satisfies the conditions of the theorem then

$$\sum_i \pi_i p_{ij} = \sum_i \pi_j p_{ji} = \pi_j \sum_i p_{ji} = \pi_j$$

and so  $\pi = \pi P$  and the distribution is stationary. □

The conditions  $\pi_i p_{ij} = \pi_j p_{ji}$  for all  $i, j \in S$  are known as the **local balance** (or **detailed balance**) conditions.



## Ehrenfest model

Suppose we have two containers  $A$  and  $B$  containing a total of  $m$  balls. At each time step a ball is chosen uniformly at random and switched between containers. Let  $X_n$  be the number of balls in container  $A$  after  $n$  units of time. Thus,  $(X_n)$  is a MC with transition matrix given by

$$p_{i,i+1} = 1 - \frac{i}{m}, \quad p_{i,i-1} = \frac{i}{m}.$$

Instead of solving the equations  $\pi = \pi P$  we look for solutions to

$$\pi_i p_{ij} = \pi_j p_{ji}$$

which yields  $\pi_i = \binom{m}{i} \left(\frac{1}{2}\right)^m$ , a binomial distribution with parameters  $m$  and  $\frac{1}{2}$ .

## Random walk on an undirected graph

Consider a **graph**  $G$  consisting of a countable collection of vertices  $i \in N$  and a finite collection of edges  $(i, j) \in E$  joining (unordered) pairs of vertices. Assume also that  $G$  is connected. A natural way to construct a MC on  $G$  uses a random walk through the vertices. Let  $v_i$  be the number of edges incident at vertex  $i$ . The random walk then moves from vertex  $i$  by selecting one of the  $v_i$  edges with equal probability  $1/v_i$ . So the transition matrix,  $P$ , is

$$p_{ij} = \begin{cases} \frac{1}{v_i} & \text{if } (i, j) \text{ is an edge} \\ 0 & \text{otherwise.} \end{cases}$$

Since  $G$  is connected,  $P$  is irreducible. The local balance conditions for  $(i,j) \in E$  are

$$\pi_i p_{ij} = \pi_j p_{ji}$$

$$\pi_i \frac{1}{v_i} = \pi_j \frac{1}{v_j}$$

$$\frac{\pi_i}{\pi_j} = \frac{v_j}{v_i}.$$

Hence,

$$\pi_i \propto v_i$$

and the normalization condition  $\sum_{i \in N} \pi_i = 1$  gives

$$\pi_i = \frac{v_i}{\sum_{j \in N} v_j}$$

and  $P$  is reversible.

# Ergodic results

Ergodic results tell us about the limiting behaviour of averages taken over time. In the case of Markov Chains we shall consider the long-run proportion of time spent in a given state.

Let  $V_i(n)$  be the **number of visits to  $i$  before time  $n$**  then

$$V_i(n) = \sum_{k=0}^{n-1} \mathbb{I}(\{X_k = i\}).$$

Thus,  $V_i(n)/n$  is the **proportion of time spent in state  $i$  before time  $n$** .

## Theorem (Ergodic theorem)

*Let  $(X_n)$  be a MC with irreducible transition matrix  $P$  then*

$$\mathbb{P}\left(\frac{V_i(n)}{n} \rightarrow \frac{1}{\mu_i} \quad \text{as } n \rightarrow \infty\right) = 1$$

*where  $\mu_i = \mathbb{E}(T_i | X_0 = i)$  is the expected return time to state  $i$ .*

## Proof

If  $P$  is transient then the total number of visits,  $V_i$ , to  $i$  is finite with probability one, so

$$\frac{V_i(n)}{n} \leq \frac{V_i}{n} \rightarrow 0 = \frac{1}{\mu_i} \quad n \rightarrow \infty.$$

Alternatively, if  $P$  is recurrent let  $Y_i^{(r)}$  be the  $r^{\text{th}}$  duration between visits to any given state  $i$ . Then  $Y_i^{(1)}, Y_i^{(2)}, \dots$  are non-negative IID RVs with  $\mathbb{E}(Y_i^{(r)}) = \mu_i$ .

But

$$Y_i^{(1)} + \dots + Y_i^{(V_i(n)-1)} \leq n - 1$$

since the time of the last visit to  $i$  before time  $n$  occurs no later than time  $n - 1$  and

$$Y_i^{(1)} + \dots + Y_i^{(V_i(n))} \geq n$$

since the time of the first visit to  $i$  after time  $n - 1$  occurs no earlier than time  $n$ .

Hence,

$$\frac{Y_i^{(1)} + \dots + Y_i^{(V_i(n)-1)}}{V_i(n)} \leq \frac{n}{V_i(n)} \leq \frac{Y_i^{(1)} + \dots + Y_i^{(V_i(n))}}{V_i(n)}.$$

However, by the SLLN,

$$\mathbb{P}\left(\frac{Y_i^{(1)} + \dots + Y_i^{(n)}}{n} \rightarrow \mu_i \text{ as } n \rightarrow \infty\right) = 1$$

and for  $P$  recurrent we know that  $\mathbb{P}(V_i(n) \rightarrow \infty \text{ as } n \rightarrow \infty) = 1$ . So,

$$\mathbb{P}\left(\frac{n}{V_i(n)} \rightarrow \mu_i \text{ as } n \rightarrow \infty\right) = 1$$

which implies

$$\mathbb{P}\left(\frac{V_i(n)}{n} \rightarrow \frac{1}{\mu_i} \text{ as } n \rightarrow \infty\right) = 1.$$



## Example: random surfing on web graphs

Consider a web graph,  $G = (V, E)$ , with vertices given by a finite collection of web pages  $i \in V$  and (directed) edges given by  $(i, j)$  whenever there is a hyperlink from page  $i$  to page  $j$ .

Random walks through the web graph have received much attention in the last few years.

Consider the following model, let  $X_n \in V$  be the location (that is, web page visited) by the surfer at time  $n$  and suppose we choose  $X_{n+1}$  uniformly from the,  $L(i)$ , outgoing links from  $i$ , in the case where  $L(i) > 0$  and uniformly among all pages in  $V$  if  $L(i) = 0$  (the **dangling page** case).

Hence, the transition matrix,  $\hat{P}_{ij}$ , say, is given by

$$\hat{p}_{ij} = \begin{cases} \frac{1}{L(i)} & \text{if } (i,j) \in E \\ \frac{1}{|V|} & \text{if } L(i) = 0 \\ 0 & \text{otherwise} \end{cases}$$

where  $|V|$  is the number of pages (that is, vertices) in the web graph. A potential problem remains in that  $\hat{P}$  may not be irreducible or may be periodic.



We will make a further adjustment to ensure irreducibility and aperiodicity as follows. For  $0 < \alpha \leq 1$  set

$$p_{ij} = (1 - \alpha)\hat{p}_{ij} + \alpha \frac{1}{|V|}.$$

We can interpret this as an “easily bored web surfer” model and see that the transitions take the form of a mixture of two distributions. With probability  $1 - \alpha$  we follow the randomly chosen outgoing link (unless the page is dangling in which case we move to a randomly chosen page) while with probability  $\alpha$  we jump to a random page selected uniformly from the entire set of pages  $V$ .

# PageRank

Brin *et al* (1999) used this approach to define PageRank through the limiting distribution of this Markov Chain, that is  $\pi_i$  where the vector  $\pi$  satisfies

$$\pi = \pi P$$

They report typical values for  $\alpha$  of between 0.1 and 0.2.

The ergodic theorem now tells us that the random surfer in this model spends a proportion  $\pi_i$  of the time visiting page  $i$  — a notion in some sense of the **importance** of page  $i$ .

Thus, two pages  $i$  and  $j$  can be ranked according to the total order defined by

$$i \geq j \quad \text{if and only if} \quad \pi_i \geq \pi_j.$$

See, “The PageRank Citation Ranking: Bring Order to the Web” Sergey Brin, Lawrence Page, Rajeev Motwani and Terry Winograd (1999) Technical Report, Computer Science Department, Stanford University.

<http://dbpubs.stanford.edu:8090/pub/1999-66>

# Computing PageRank: the power method

We seek a solution to the system of equations

$$\pi = \pi P$$

that is, we are looking for an eigenvector of  $P$  (with corresponding eigenvalue of one). Google's computation of PageRank is one of the world's largest matrix computations.

The power method starts from some initial distribution  $\pi^{(0)}$ , updating  $\pi^{(k-1)}$  by the iteration

$$\pi^{(k)} = \pi^{(k-1)} P = \dots = \pi^{(0)} P^k$$

Advanced methods from linear algebra can be used to speed up convergence of the power method and there has been much study of related MCs to include web browser back buttons and many other properties as well as alternative notions of the “importance” of a web page.

# Hidden Markov Models

An extension of Markov Chains is provided by **Hidden Markov Models** (HMM) where a statistical model of observed data is constructed from an underlying but usually hidden Markov Chain.

Such models have proved very popular in a wide variety of fields including

- ▶ speech and optical character recognition
- ▶ natural language processing
- ▶ bioinformatics and genomics.

We shall not consider these applications in any detail but simply introduce the basic ideas and questions that Hidden Markov Models address.

## A Markov model with hidden states

Suppose we have a MC with transition matrix  $P$  but that the states  $i$  of the chain are not directly observable. Instead, we suppose that on visiting any state  $i$  at time  $n$  there is a randomly chosen output value or token,  $Y_n$ , that is observable.

The probability of observing the output token  $t$  when in state  $i$  is given by some distribution  $b_i$ , depending on the state  $i$  that is visited.

Thus,

$$\mathbb{P}(Y_n = t | X_n = i) = (b_i)_t$$

where  $(b_i)_t$  is the  $t^{\text{th}}$  component of the distribution  $b_i$ .

For an excellent introduction to HMM, see “A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition” Lawrence R. Rabiner. Proceedings of the IEEE, Vol 77, No 2, February 1988.

# Three central questions

There are many variants of this basic setup but three central problems are usually addressed.

## Definition (Evaluation problem)

Given a sequence  $y_1, y_2, \dots, y_n$  of observed output tokens and the parameters of the HMM (namely,  $P$ ,  $b_i$  and the distribution for the initial state  $X_0$ ) how do we compute

$$\mathbb{P}(Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n | \text{HMM parameters})$$

that is, the probability of the observed sequence given the model?

Such problems are solved in practice by the **forward algorithm**.

A second problem that may occur in an application is the **decoding problem**.

### **Definition (Decoding problem)**

Given an observed sequence of output tokens  $y_1, y_2, \dots, y_n$  and the full description of the HMM parameters, how do we find the best fitting corresponding sequence of (hidden) states  $i_1, i_2, \dots, i_n$  of the MC?

Such problems are solved in practice by a dynamic programming approach called the **Viterbi algorithm**.

The third important problem is the **learning problem**.

### Definition (Learning problem)

Given an observed sequence of output tokens  $y_1, y_2, \dots, y_n$ , how do we adjust the parameters of the HMM to maximize

$$\mathbb{P}(Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n | \text{HMM parameters})$$

The observed sequence used to adjust the model parameters is called a **training sequence**. Learning problems are crucial in most applications since they allow us to create the “**best**” models in real observed processes.

Iterative procedures, known as the **Baum-Welch method**, are used to solve this problem in practice.



# Applications of Markov Chains

These and other applications of Markov Chains are important topics in a variety of Part II courses, including

- ▶ Artificial Intelligence II
- ▶ Bioinformatics
- ▶ Computer Systems Modelling

# Case studies

# Case studies

Three short cases studies where probability has played a pivotal role:

1. Birthday problem (**birthday attack**)
  - cryptographic attacks
2. Probabilistic classification (**naive Bayes classifier**)
  - email spam filtering
3. Gambler's ruin problem (**Bitcoin**)
  - cryptocurrencies

# The birthday problem

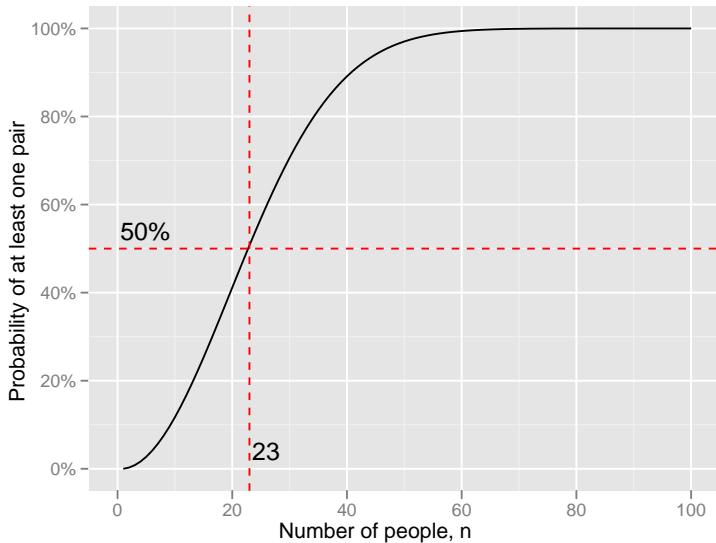
Consider the problem of computing the probability,  $p(n)$ , that in a party of  $n$  people at least two people share a birthday (that is, the same day and month but not necessarily same year).

It is easiest to first work out  $1 - p(n) = q(n)$ , say, where  $q(n) = \mathbb{P}(\text{none of the } n \text{ people share a birthday})$  then

$$\begin{aligned} q(n) &= \left(\frac{364}{365}\right) \left(\frac{363}{365}\right) \cdots \left(\frac{365-n+1}{365}\right) \\ &= \left(1 - \frac{1}{365}\right) \left(1 - \frac{2}{365}\right) \cdots \left(1 - \frac{n-1}{365}\right) \\ &= \prod_{k=1}^{n-1} \left(1 - \frac{k}{365}\right). \end{aligned}$$

Surprisingly,  $n = 23$  people suffice to make  $p(n)$  greater than 50%.

# Graph of $p(n)$



# Assumptions

We should record some of our assumptions behind the calculation of  $p(n)$ .

1. Ignore leap days (29 Feb)
2. Each birthday is equally likely
3. People are selected independently and without regard to their birthday to attend the party (ignore twins, etc)

## Examples: coincidences on the football field

Ian Stewart writing in Scientific American illustrates the birthday problem with an interesting example. In a football match there are 23 people (two teams of 11 plus the referee) and on 19 April 1997 out of 10 UK Premier Division games there were 6 games with birthday coincidences and 4 games without.

# Examples: cryptographic hash functions

A hash function  $y = f(x)$  used in cryptographic applications is usually required to have the following two properties (amongst others):

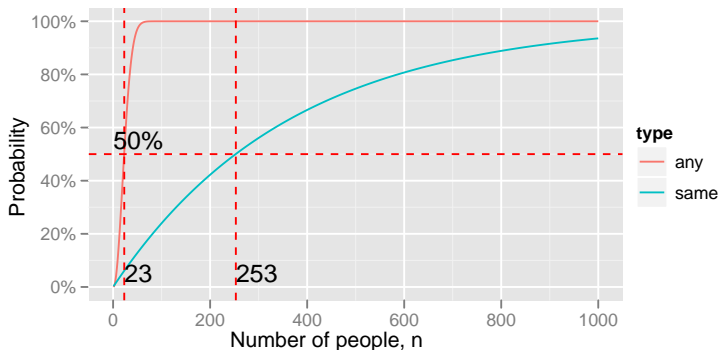
1. **one-way function**: computationally intractable to find an  $x$  given  $y$ .
2. **collision-resistant**: computationally intractable to find distinct  $x_1$  and  $x_2$  such that  $f(x_1) = f(x_2)$ .



# Probability of same birthday as you

Note that in calculating  $p(n)$  we are not specifying which birthday (for example, your own) matches. For the case of finding a match to your own birthday amongst a party of  $n$  other people we would calculate

$$1 - \left(\frac{364}{365}\right)^n.$$



## General birthday problem

Suppose we have a random sample  $X_1, X_2, \dots, X_n$  of size  $n$  where  $X_i$  are IID with  $X_i \sim U(1, d)$  and let  $p(n, d)$  be the probability that there are at least two outcomes that coincide.

Then

$$p(n, d) = \begin{cases} 1 - \prod_{k=1}^{n-1} \left(1 - \frac{k}{d}\right) & n \leq d \\ 1 & n > d. \end{cases}$$

The usual birthday problem is the special case when  $d = 365$ .

# Approximations

One useful approximation is to note that for  $x \ll 1$  then  $1 - x \approx e^{-x}$ .  
Hence for  $n \leq d$

$$\begin{aligned} p(n, d) &= 1 - \prod_{k=1}^{n-1} \left(1 - \frac{k}{d}\right) \\ &\approx 1 - \prod_{k=1}^{n-1} e^{-\frac{k}{d}} \\ &= 1 - e^{-(\sum_{k=1}^{n-1} k)/d} \\ &= 1 - e^{-n(n-1)/(2d)}. \end{aligned}$$

We can further approximate the last expression as

$$p(n, d) \approx 1 - e^{-n^2/(2d)}.$$

# Inverse birthday problem

Using the last approximation

$$p(n, d) \approx 1 - e^{-n^2/(2d)}$$

we can invert the birthday problem to find  $n = n(p, d)$ , say, such that  $p(n, d) \approx p$  so then

$$e^{-n(p, d)^2/(2d)} \approx 1 - p$$

$$-\frac{n(p, d)^2}{2d} \approx \log(1 - p)$$

$$n(p, d)^2 \approx 2d \log\left(\frac{1}{1 - p}\right)$$

$$n(p, d) \approx \sqrt{2d \log\left(\frac{1}{1 - p}\right)}.$$

In the special case of  $d = 365$  and  $p = 1/2$  this gives the approximation  $n(0.5, 365) \approx \sqrt{2 \times 365 \times \log(2)} \approx 22.49$ .

## Expected waiting times for a collision/match

Let  $W_d$  be the random variable specifying the number of iterations when you choose one of  $d$  values independently and uniformly at random (with replacement) and stop when any value is selected a second time (that is, a “collision” or “match” occurs).

It is possible to show that

$$\mathbb{E}(W_d) \approx \sqrt{\frac{\pi d}{2}}.$$

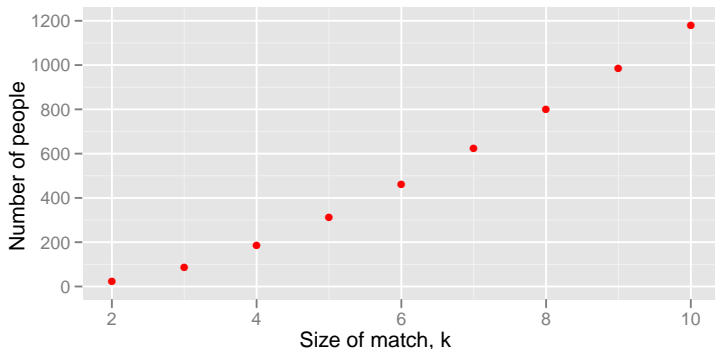
Thus in the special case of the birthday problem where  $d = 365$  we have that  $\mathbb{E}(W_{365}) \approx \sqrt{\frac{\pi \times 365}{2}} \approx 23.94$ .

In the case that we have a cryptographic hash function with 160-bit outputs ( $d = 2^{160}$ ) then  $\mathbb{E}(W_{2^{160}}) \approx 1.25 \times 2^{80}$ . This level of reduction leads to so-called “**birthday attacks**”. (See the IB course Security I for further details.)

## Further results

Persi Diaconis and Frederick Mosteller give results on the minimum number  $n_k$  required to give a probability greater than  $1/2$  of  $k$  or more matches with  $d = 365$  possible choices.

$k$	2	3	4	5	6	7	8	9	10
$n_k$	23	88	187	313	460	623	798	985	1181



## Email spam filtering

Suppose that an email falls into exactly one of two classes (spam or ham) and that various features  $F_1, F_2, \dots, F_n$  of an email message can be measured. Such features could be the presence or absence of particular words or groups of words, etc, etc.

We would like to determine  $\mathbb{P}(C | F_1, F_2, \dots, F_n)$  the probability that an email message falls into a class  $C$  given the measured features  $F_1, F_2, \dots, F_n$ . We can use Bayes' theorem to help us.

# Bayes' theorem for emails

We have that

$$\mathbb{P}(C|F_1, F_2, \dots, F_n) = \frac{\mathbb{P}(C)\mathbb{P}(F_1, F_2, \dots, F_n|C)}{\mathbb{P}(F_1, F_2, \dots, F_n)}$$

which can be expressed in words as

$$\text{posterior probability} = \frac{\text{prior probability} \times \text{likelihood}}{\text{evidence}}.$$



# Naive Bayes classifier

In the **naive Bayes classifier** we make the assumption of independence across features. So that

$$\mathbb{P}(F_1, F_2, \dots, F_n | C) = \prod_{i=1}^n \mathbb{P}(F_i | C)$$

and then

$$\mathbb{P}(C | F_1, F_2, \dots, F_n) \propto \mathbb{P}(C) \prod_{i=1}^n \mathbb{P}(F_i | C).$$

## Decision rule for naive Bayes classifier

We then use the **decision rule** to classify an email with observed features  $F_1, F_2, \dots, F_n$  as spam if

$$\mathbb{P}(C = \text{spam}) \prod_{i=1}^n \mathbb{P}(F_i | C = \text{spam}) > \mathbb{P}(C = \text{ham}) \prod_{i=1}^n \mathbb{P}(F_i | C = \text{ham}).$$

This decision rule is known as the **maximum a posteriori** (MAP) rule. Surveys and a training set of manually classified emails are needed to estimate the values of  $\mathbb{P}(C)$  and  $\mathbb{P}(F_i | C)$ .

# Bitcoin

Bitcoin is based around a **proof-of-work** mechanism which uses a decentralised peer-to-peer network of workers (known as **miners**) to ensure (with high probability) that bitcoins are not **double-spent**. In order to achieve double spending of a bitcoin the attacker would need to create a longer block chain than the honest chain.

Suppose that the honest workers can produce blocks on average every  $T/p$  time units while the attacker can do so on average every  $T/q$  time units with  $q = 1 - p < p$ . If the (honest) seller waits for a given number  $n$  of blocks to be created then this would take on average  $nT/p$  time units. Thus the average number of blocks that the attacker could create,  $m$ , would be such that  $nT/p = mT/q$ .

Thus  $m = nq/p$  independent of  $T$ .

If  $q > p$  then surely the attacker can always catch up the honest workers however large a head start,  $n$ , is considered. What is the chance that the attacker could still catch up the honest chain when  $q < p$ ?

## Bitcoin analysis using Gambler's ruin problem

The Bitcoin white paper proposes the simple probabilistic model that the random number of blocks,  $X$ , that the attacker could produce as the honest workers produce their  $n$  blocks has a Poisson distribution with mean  $\lambda = nq/p$ . Thus,

$$\mathbb{P}(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}.$$

What is the chance that the attacker could then overtake the honest workers? This is precisely the Gambler's ruin problem starting from initial assets of  $n - k$  with  $\theta = q/p < 1$  and in the limit that the total wealth  $N \rightarrow \infty$ .

Recalling our expression for the ruin probabilities we have that

$$\mathbb{P}(\text{attacker catches up} \mid k \text{ blocks}) = \begin{cases} (q/p)^{n-k} & k \leq n \\ 1 & k > n. \end{cases}$$

## Bitcoin: calculations

Hence, using the law of total probability,

$$\mathbb{P}(\text{attacker catches up}) = \sum_{k=0}^{\infty} \mathbb{P}(X = k) \begin{cases} (q/p)^{n-k} & k \leq n \\ 1 & k > n \end{cases}$$

which we can re-write as the finite sum

$$1 - \sum_{k=0}^n \mathbb{P}(X = k) \left(1 - \left(\frac{q}{p}\right)^{n-k}\right) = 1 - \sum_{k=0}^n \frac{\lambda^k e^{-\lambda}}{k!} \left(1 - \left(\frac{q}{p}\right)^{n-k}\right)$$

where  $\lambda = nq/p$  is the mean number of blocks that an attacker can produce while the honest workers produce  $n$  blocks.

The bitcoin white paper (section 11) provides a C program to compute this probability for fixed  $q = 1 - p$  and varying  $n$  and observes that the probability of catching up the honest workers drops off exponentially with  $n$ .

What other probabilistic models would you suggest in place of the Poisson assumption?

# References for case studies



Ian Stewart

*What a coincidence!*

Mathematical Recreations, Scientific American, Jun 1998, 95–96.



Persi Diaconis and Frederick Mosteller

*Methods for studying coincidences.*

Journal of American Statistical Association, Vol 84, No 408, Dec 1989, 853–861.



Satoshi Nakamoto

*Bitcoin: A peer-to-peer electronic cash system.*

<http://bitcoin.org/bitcoin.pdf>, 2008.

# Properties of discrete RVs

RV, $X$	Parameters	$\text{Im}(X)$	$\mathbb{P}(X = k)$	$\mathbb{E}(X)$	$\text{Var}(X)$	$G_X(z)$
Bernoulli	$p \in [0, 1]$	$\{0, 1\}$	$(1-p)$ if $k = 0$ or $p$ if $k = 1$	$p$	$p(1-p)$	$(1-p+pz)$
$\text{Bin}(n, p)$	$n \in \{1, 2, \dots\}$ $p \in [0, 1]$	$\{0, 1, \dots, n\}$	$\binom{n}{k} p^k (1-p)^{n-k}$	$np$	$np(1-p)$	$(1-p+pz)^n$
$\text{Geo}(p)$	$0 < p \leq 1$	$\{1, 2, \dots\}$	$p(1-p)^{k-1}$	$\frac{1}{p}$	$\frac{1-p}{p^2}$	$\frac{pz}{1-(1-p)z}$
$U(1, n)$	$n \in \{1, 2, \dots\}$	$\{1, 2, \dots, n\}$	$\frac{1}{n}$	$\frac{n+1}{2}$	$\frac{n^2-1}{12}$	$\frac{z(1-z^n)}{n(1-z)}$
$\text{Pois}(\lambda)$	$\lambda > 0$	$\{0, 1, \dots\}$	$\frac{\lambda^k e^{-\lambda}}{k!}$	$\lambda$	$\lambda$	$e^{\lambda(z-1)}$

# Properties of continuous RVs

RV, $X$	Parameters	$\text{Im}(X)$	$f_X(x)$	$\mathbb{E}(X)$	$\text{Var}(X)$
$U(a, b)$	$a, b \in \mathbb{R}$ $a < b$	$(a, b)$	$\frac{1}{b-a}$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
$\text{Exp}(\lambda)$	$\lambda > 0$	$\mathbb{R}_+$	$\lambda e^{-\lambda x}$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$
$N(\mu, \sigma^2)$	$\mu \in \mathbb{R}$ $\sigma^2 > 0$	$\mathbb{R}$	$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/(2\sigma^2)}$	$\mu$	$\sigma^2$