# **Basics**

To recap, any probability space is a tuple of a sample space, a collection of subsets of that sample space, and a function from the event space to [0,1], where the event space and probability function each satisfy three natural conditions ( $\mathcal{F}$  is a  $\sigma$ -algebra, plus the basic properties of  $\mathbb{P}$ ). We take random variables as functions from  $\Omega$ , representing an observable. Formally they must also have that each  $\{X(\omega) \leq x\} \in \mathcal{F}$ .

In prelims probability there was a distinction between discrete and continuous random variables. These do not cover every possible notion of a random variable, and so we ideally want to unify these definitions to a more abstract notion. Beginning in this way we define

## Definition 1 (Expectation) • $\mathbb{E}I_A = \mathbb{P}(A)$ for any event A.

• If  $\mathbb{P}(X \ge 0) = 1$  then  $\mathbb{E}X \ge 0$ .

•  $\mathbb{E}(X + aY) = \mathbb{E}X + a\mathbb{E}Y \text{ for any } a \in \mathbb{R}.$ We immediately get consequences of these axioms for notions of variance and covariance,

so we need not add additional baggage to each of these for the moment.

**Definition 2 (Independence)** A collection of events  $\{A_i | i \in I\}$  are independent if  $\mathbb{P}\left(\bigcap_{i\in I}A_i\right)=\prod_{i\in I}\mathbb{P}(A_i)$ 

### Take X, Y random variables. In certain cases we would like a concept of distance between X and Y.

Convergence of random variables

**Definition 3 (Convergence)** Take a sequence  $(X_n)$  of random variables, and random variable X.

X<sub>n</sub> → X (almost surely) if P({X<sub>n</sub> → X as n → ∞}) = 1.
X<sub>n</sub> → X (in probability) as n → ∞ if for every ε > 0, P(|X<sub>n</sub> - X| < ε) → 1 as</li>  $n \to \infty$ .

•  $X_n \stackrel{d}{\to} X$  (in distribution) as  $n \to \infty$  if for every  $x \in \mathbb{R}$  such that F is continuous at  $x, F_n(x) \to F(x)$  as  $n \to \infty$ .

We should find that the above notions are decreasing in strength. By its nature we can often write distribution convergence not with a random variable X, but just with its distribution. To show that almost sure convergence implies probabilistic convergence, we first state the

**Lemma 1** Let  $A_n$  be an increasing sequence of events (For all  $k \in \mathbb{N}$ ,  $A_k \subseteq A_{k+1}$ ). Then $\mathbb{P}(A_n) \to \mathbb{P}\left(\bigcup_{k=0}^{\infty} A_k\right)$ 

As proof, write

 $\mathbb{P}(A_n) = \mathbb{P}\left(\bigcup_{k=0}^n A_k\right) = \mathbb{P}\left(A_0 \cup \bigcup_{k=1}^n A_k \setminus A_{k-1}\right)$ 

$$k=0 \qquad k=1$$

$$= \mathbb{P}(A_0) + \sum_{k=1}^{n}$$

following lemma:

$$= \mathbb{P}(A_0) + \sum_{k=1}^n \mathbb{P}(A_k \setminus A_{k-1})$$

$$\to \mathbb{P}(A_0) + \sum_{k=1}^\infty \mathbb{P}(A_k \setminus A_{k-1})$$

$$= \mathbb{P}\left(\bigcup_{k=0}^\infty A_k\right).$$
We can then consider the event defined in almost sure convergence:
$$\{X_n \to X \text{ as } n \to \infty\} = \{\forall \varepsilon > 0 . \exists N \ge 0 . \forall n \ge N . |X_n - X| < \varepsilon\}$$

$$= \bigcap_{k=0}^\infty \{\forall n \ge N . |X_n - X| < \varepsilon\}$$

 $\subseteq \bigcup \{ \forall n \ge N . |X_n - X| < \varepsilon \}$  for any  $\varepsilon > 0$ Thus we turn the event of convergence into an infinite union of increasing sets, which is

tail being within a small range being likely. See  $X_n \sim \text{Ber}(1/n)$ .

probabilistic convergence is achieved.

constant c implies that  $X_n \stackrel{p}{\to} c$ .

Martingales.

itself an event of probability 1, so we have 
$$\mathbb{P}(\{\forall n \geq N : |X_n - X| < \varepsilon\}) \to 1$$
 as  $N \to \infty$ . Further, 
$$\{\forall n \geq N : |X_n - X| < \varepsilon\} = \bigcap_{n=N}^{\infty} \{|X_n - X| < \varepsilon\}$$
  $\subseteq \{|X_n - X| < \varepsilon\}$  for any  $n \geq N$  so we get  $1 \geq \mathbb{P}(|X_n - X| < \varepsilon) \geq \mathbb{P}(\{\forall n \geq N : |X_n - X|\}) \to 1$  and by sandwiching

To show that probabilistic convergence implies distributive convergence, note that in the limit we can get  $F_n(x)$  in terms of an arbitrary  $\varepsilon > 0$  and X. Then we may bound  $F_n(x)$ and use continuity of F to show convergence.

**Theorem 2** For  $(X_n)$  all defined on the same probability space,  $X_n \stackrel{d}{\rightarrow} c$  for some

To show that the inverse doesn't hold, just take a sequence of random variables wherein

the probability clearly converges, but not so quickly as to have the probability of an infinite

This follows fairly immediately from algebra. **Theorem 3 (Weak law of large numbers)** Suppose  $(X_n)$  are i.i.d. with mean  $\mu < \infty$ . Let  $S_n = \sum_{k=1}^n X_k$ . Then

 $\frac{S_n}{m} \xrightarrow{p} \mu \ as \ n \to \infty$ 

 $\phi_{S_n/n}(t) = \phi_X(t/n)^n$ 

We can prove this statement using characteristic functions:

 $= \left(1 + i\mathbb{E}[X]\frac{t}{n} + o(t/n)\right)^n$  $\rightarrow e^{it\mathbb{E}[X]}$  by continuity of exp and log

and by the characteristic function continuity result  $S_n/n \stackrel{d}{\to} \mu$ , which then means

f large numbers) 
$$Suppose(X_n)$$
 are  $ii$ 

 $S_n/n \xrightarrow{p} \mu$  as  $\mu$  is constant. **Theorem 4 (Strong law of large numbers)** Suppose  $(X_n)$  are iid with mean  $\mu < \infty$ . Let  $S_n = \sum_{k=1}^n X_k$ . Then

$$\frac{S_n}{n} \to \mu \ almost \ surely \ as \ n \to \infty$$
 The proof of this is not examinable, and a full proof is given in Probability, Measure and

 $\sigma^2 < \infty$ . Let  $S_n = \sum_{k=1}^n X_k$ , then  $\frac{S_n - n\mu}{\sigma \sqrt{n}} \xrightarrow{d} N(0,1) \text{ as } n \to \infty$ 

**Theorem 5 (Central limit theorem)** Suppose  $(X_n)$  are i.i.d.,  $\mathbb{E}[X_k] = \mu$ ,  $\operatorname{Var} X_k = 0$ 

 $\phi_{S_n/\sqrt{n}}(t) = \phi_Y(\frac{t}{\sqrt{n}})^n$  $= \left(1 - \frac{t^2}{2n} + o(t^2/n)\right)^n$   $\to e^{-t^2/2}$ 

so by continuity 
$$S_n/\sqrt{n} \stackrel{d}{\to} N(0,1).$$

Conditional Densities

Definition 4 For two events A and B with  $\mathbb{P}(A) > 0$ ,

We define  $Y_n = \frac{X_n - \mu}{\sigma}$ , so  $S_n = \frac{1}{\sqrt{n}} \sum_{k=1}^n Y_k$ , and thus

 $\mathbb{P}(B \mid A) := \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(A)},$ and in application to random variables, we get  $\mathbb{P}(X \le x \mid A) = \frac{\mathbb{P}(\{X \le x\} \cap A)}{\mathbb{P}(A)}$ 

## The second function gives a conditional cdf for X, implying the existence of a pdf $f_{X|A}$ for which $\mathbb{P}(X \in C \mid A) = \int_C f_{X|A}(x) \, \mathrm{d}x$

Markov Chains

distribution of  $X_{n+1}$  given  $X_n = i$ .

of Y conditioned on  $\{x \leq X \leq x + \varepsilon\}$ , and for nice enough  $f_{X,Y}(x,y)$ ,  $f_X(x)$  we get  $\mathbb{P}(Y \le y \mid x \le X \le x + \varepsilon) = \frac{\int_{-\infty}^{y} \int_{x}^{x+\varepsilon} f_{X,Y}(u, v) \, \mathrm{d}u \, \mathrm{d}v}{\int_{x}^{x+\varepsilon} f_{X}(u) \, \mathrm{d}u}$  $\sim \int_{-\infty}^{y} \frac{f_{X,Y}(x, v)}{f_{X}(x)} \, \mathrm{d}v \quad \text{as } \varepsilon \to 0$ 

in I. The process X is called a Markov chain if for any 
$$n \geq 0$$
 and  $i_0, i_1, \ldots, i_{n+1} \in I$ , 
$$\mathbb{P}(X_{n+1} = i_{n+1} \mid X_n = i_n, \ldots, X_0 = i_0) = \mathbb{P}(X_{n+1} = i_{n+1} \mid X_n = i_n)$$
In addition, the Markov chain is homogeneous if  $\mathbb{P}(X_{n+1} = j \mid X_n = i)$  is constant in  $n \geq 0$ .
Intuitively, a Markov chain is a sequence wherein one need not keep track of previous

**Definition 5** Let  $X = (X_0, X_1, X_2, \dots)$  be a sequence of random variables taking values

states in order to determine the distribution over future states, but rather one only needs

to know where they are (and potentially the time at which they are there). In the case of

a homogeneous Markov chain, we can write  $P = (p_{ij})$  as the matrix with the ith row the

We almost always talk about homogeneous Markov chains in this course.

A problem which we come to is trying to observe the conditional density of Y for X=x,

as for continuous random variables  $\mathbb{P}(X=x)=0$ . To resolve this, we take the distribution

**Theorem 6 (\*Extended Markov Property)** Let  $(X_n)$  be a Markov chain. For  $n \geq 0$ , for any event H given in terms of the past history  $X_0, X_1, \ldots, X_{n-1}$ , and any event F given in terms of  $X_{n+1}, X_{n+2}, \ldots$ , we have  $\mathbb{P}(F \mid X_n = i, H) = \mathbb{P}(F \mid X_n = i)$ To prove this for F a statement about an infinite future requires material beyond this

course, but for finite cases we can just write out the conditional to get the result.

From the Markov property, we very quickly get a formula for n-step probabilities.

 $p_{ij}^{(n+m)} = \mathbb{P}(X_{n+m+r} = j \mid X_r = i)$   $= \sum_{k \in I} \mathbb{P}(X_{m+r} = k \mid X_r = i) \mathbb{P}(X_{n+m+r} = j \mid X_{m+r} = k)$ 

It is not quite correct to say that in a Markov chain  $X_n$  depends only on  $X_{n-1}$  - there is

certainly still randomness involved, and this would imply a functional relationship which

doesn't quite exist. We can however say that for each n we can have a random variable

 $Y_n = f(Y_{n-1}, X_n)$  where  $X_n$  is independent of  $(Y_0, \dots, Y_{n+1})$ . Then  $(Y_n)$  is a markov chain.

We say that i leads to j, or  $i \to j$  where for some  $n \ge 0$ ,  $p_{ij}^{(n)} > 0$ , and we say that i

communicates with j, or  $i \leftrightarrow j$  where  $i \to j$  and  $j \to i$ . This is an equivalence relation,

thus partitioning I into communicating classes. We say that a chain for which I is a single

equivalence class is irreducible. Further we say that a class is closed if the probability for

**Definition 6 (Period)** The periodicity of state i is defined as  $gcd\{n \mid p_{ii}^{(n)} > 0\}$ . If this

ever exiting is 0. If the singleton of a state is closed then that state is absorbing.

 $=\sum_{k\in I}p_{ik}^{(m)}p_{kj}^{(n)}$  $= (P^{(m)}P^{(n)})_{ij}$ so  $P^{(n)} = P^{(n-1)}P$  so by induction  $P^{(n)} = P^n$ .

All states within the same communicating class have the same period. To see this note that if 
$$i$$
 and  $j$  communicate, then we can get  $a$ ,  $b$  such that  $p_{ij}^{(a)} > 0$  and  $p_{ji}^{(b)} > 0$ , so  $p_{ii}^{(a+b)} > 0$ . Further, if  $p_{jj}^{(m)} > 0$ , then  $p_{ii}^{(a+b+m)} > 0$ . Thus if  $i$  has period  $d$ , then  $d \mid a+b+m$  and  $d \mid a+b$ , so  $d \mid m$ . Thus the period of  $i$  divides the period of  $j$ , and by symmetry thus the reverse holds, so the period of  $i$  is equal to the period of  $j$ .

Definition 7 Let  $(X_n)$  be a Markov chain, and  $A \subseteq I$ . Define

 $h_i^A = \mathbb{P}\left(\bigcup_{n\geq 0} \{X_n \in A\} \mid X_0 = i\right)$ 

as the hitting probability of A from i. **Theorem 7** The vector of hitting probabilities  $(h_i^A | i \in I)$  is the minimal non-negative

solution  $\boldsymbol{x}$ , and show that for all  $M \in \mathbb{N}$ ,  $i \in I$  that

solution to the recurrence equations

is 1, so we call the state recurrent.

recurrent.

$$h_i^A = \begin{cases} 1 & \text{if } i \in A \\ \sum_{j \in I} p_{ij} h_j^A & \text{if } i \notin A \end{cases}$$
 The base case is obvious. For the recurrence we partition and use the Markov property. To show that the minimal non-negative solution is correct, take an arbitrary non-negative

For M=0 we get if  $i\in A$  that  $x_i=1$ , and if  $i\notin A$  that the right hand side is 0. Further, if the statement is true for M-1, then if  $i \in A$  then again  $x_i = 1$  so the equation holds, and otherwise we can partition to maintain the inequality. Recurrence and Transience

For  $\mathbb{P}(X_n = i \text{ for some } n \geq 1) < 1$ , we have that the total number of visits to

i has geometric distribution with parameter 1-p, and so the probability that

i is hit infinitely often is 0, so we call the state transient. If however we have

 $\mathbb{P}(X_n = i \text{ for some } n \geq 1) = 1$ , then clearly the probability of hitting i infinitely often

 $x_i \ge \mathbb{P}\left(\bigcup_{n \le M} \{X_n \in A\} \mid X_0 = i\right).$ 

The total number of visits to i is  $\sum \mathbb{1}(X_n = i)$ , which has expectation equal to  $\sum p_{ii}^{(n)}$ . If i is transient this expectation is finite, whereas if it is recurrent then the expectation

positive, so  $p_{ji}^{(a)} p_{ii}^{(n)} p_{ij}^{(b)} \leq p_{jj}^{(a+b+n)}$ , so  $\frac{1}{p_{ii}^{(a)} p_{ij}^{(b)}} \sum_{n=0}^{\infty} p_{jj}^{(n)}$  is infinite.

**Theorem 8** A state i is recurrent iff  $\sum_{n=0}^{\infty} p_{ii}^{(n)} = \infty$ .

is infinite. **Theorem 9** Let C be a communicating class. Either all states in C are recurrent, or all are transient. Further, every recurrent class is closed, and every finite closed class is recurrent. Take a C with a recurrent state, so  $\sum_{n=0}^{\infty} p_{ii}^{(n)}$  is infinite. For some  $a, b, p_{ii}^{(a)}, p_{ij}^{(b)}$  are

would return to their original position eventually with probability 1. If, however, they have access to a spaceship, then there is positive probability that they never come home.

**Definition 8**  $H^A = \min\{n \geq 0 \mid X_n \in A\}$  is the hitting time of A. **Theorem 10** The vector of mean hitting times  $k^A$  is the minimal non-negative solution to

 $k_i^A = \begin{cases} 0 & \text{if } i \in A \\ 1 + \sum_{j} p_{ij} k_j^A & \text{otherwise} \end{cases}$ The proof here follows straightforwardly from conditional expectations, and minimality

expectation axiomatically:

# Generating Functions

 $t \in [-t_0, t_0]$ :

result follows.

We have an existing notion of generating functions for discrete random variables from prelims probability. That is,  $G_X(s) = \mathbb{E}[s^X]$ , defined on the radius of convergence of the corresponding power series. We have various results about these functions, such as that the exact distribution of X may be extracted via differentiation, demonstrating uniqueness, and that with  $(X_n)$ , N independent, each  $X_n$  identically distributed,  $G_{\sum_{i=1}^N X_i}(s) = G_N(G_X(s))$ . **Theorem 11** If each  $X_n$  for  $n \geq 1$  and X have generating functions  $G_{X_n}$  and  $G_{X_n}$ 

then  $G_{X_n} \to G_X$  pointwise if and only if  $X_n \stackrel{d}{\to} X$ . This is hopefully clear from definitions.

**Definition 9** The moment generating function of a random variable X is defined as

 $M_X(t) = \mathbb{E}[e^{tX}]$ For example, for  $\text{Exp}(\lambda)$ :

example, for 
$$\mathrm{Exp}(\lambda)$$
: 
$$M_X(t) = \mathbb{E}[e^{tX}]$$

 $= \int_{-\infty}^{\infty} e^{tx} f(x) \, \mathrm{d}x$ 

$$= \int_0^\infty \lambda e^{(t-\lambda)x} \,\mathrm{d}x$$
 
$$= \begin{cases} \frac{\lambda}{\lambda - t} & \text{if } t < \lambda \\ \infty & \text{otherwise} \end{cases}$$
 We get fairly quickly a few similar results as for generating functions. For  $X$  with a generating function  $M_X$  defined for  $t$ ,

 $M_{aX+b}(t) = \mathbb{E}[e^{t(aX+b)}]$  $=e^{bt}\mathbb{E}[e^{atX}]$  $=e^{bt}M_X(at)$ 

$$m_{aX+b}(t) = \mathbb{E}[e^{tX}]$$

$$= e^{bt}\mathbb{E}[e^{atX}]$$

$$= e^{bt}M_X(at)$$
and for  $\{X_1, \dots, X_n\}$  independent with generating functions defined for each on  $t$ ,
$$M_{\sum_{k=1}^n X_k}(t) = \mathbb{E}[e^{t\sum_{k=1}^n X_k}]$$

 $= \mathbb{E}\left[\prod_{k=1}^n e^{tX_k}\right]$ 

$$=\prod_{k=1}^n\mathbb{E}[e^{tX_k}]$$
 
$$=\prod_{k=1}^nM_{X_k}(t).$$
 Furthermore, we have a convergence result, that if  $M_{|X|}(t_0)$  exists for some  $t_0>0$ , then for  $t\in[-t_0,t_0]$ : 
$$M_{|X|}(t_0)=\int_0^\infty e^{t_0x}(f(x)+f(-x))\,\mathrm{d}x$$

 $\geq \int_{0}^{\infty} e^{tx} (f(x) + f(-x)) dx \qquad \text{for } |t| \leq t_0$ 

$$\geq \int_{-\infty}^{\infty} e^{tx} f(x) dx$$

$$= M_X(t)$$
so  $M_X(t)$  is defined on this interval.

**Theorem 12** Suppose  $\mathbb{E}[e^{t_0|X|}]$  is finite for some  $t_0 > 0$ . Then we both have that
$$M_X(t) = \sum_{k=0}^{\infty} \mathbb{E}[X^k] \frac{t^k}{k!} \quad \text{for } |t| \leq t_0$$

 $M_X^{(k)}(0) = \mathbb{E}[X^k]$ 

One needs a bit of work not included in this course (Fubini's theorem) to show that the

expectation operator and infinite sums can commute in this case, but assuming that the

we can use  $\mathbb{E}[e^{tX}] \leq \mathbb{E}[e^{t|X|}]$ , from which we get  $\mathbb{E}[e^{t|X|}] = \int_0^\infty \mathbb{P}(e^{t|X|} > x) \, \mathrm{d}x$  $\leq 1 + \int_{1}^{\infty} \mathbb{P}\left(|X| > \frac{\log x}{t}\right) dx$ 

An equivalent statement to the existence of the MGF on some neighbourhood of 0 is

that for some  $t_0 > 0$ ,  $\mathbb{P}(|X| > x) = O(e^{-t_0x})$ . If  $M_X(t)$  is finite on  $[-t_0, t_0]$ , then

 $\mathbb{P}(|X|>x)\leq e^{-t_0x}M_X(t_0)$  for all  $x\geq 0$  by Markov's inequality. In the reverse direction,

$$\leq 1 + \int_{1}^{\infty} Cx^{-t_0/t} \, \mathrm{d}x$$
 which is a finite integral for  $0 < t < t_0$ .

**Theorem 13** If X and Y have the same moment generating function, which is finite on  $[-t_0, t_0]$  for some  $t_0 > 0$ , then X and Y have the same distribution.

More generally, if we have a sequence of random variables  $(X_n)$  and X with finite

moment generating functions on  $[-t_0, t_0]$ , and  $M_{X_n}(t) \to M_X(t)$  as  $n \to \infty$  for all  $t \in [-t_0, t_0], then X_n \stackrel{d}{\rightarrow} X as n \rightarrow \infty.$ 

The proofs of both the above are beyond the scope of this course. **Definition 10** The characteristic function of X is  $\phi_X(t) = \mathbb{E}[e^{itX}] = \mathbb{E}\cos(tX) + 1$  $i\mathbb{E}\sin(tX)$ .

our convergence result becomes that the characteristic function always exists. This follows as  $\cos(tX)$  and  $\sin(tX)$  have image [-1,1], so the function is just the sum of two finite integrals. Thanks to this convergence result we get the following power series result:

 $\phi_X(t) = \sum_{n=0}^{\infty} \frac{i^n t^n \mathbb{E}[X^n]}{n!}$ 

Both the uniqueness and continuity statements hold in a similar way as for MGFs, but as

Not only can we extend all of the basic results for MGFs to characteristic functions, but

before their proofs are beyond the scope of this course.

**Definition 11** The joint cumulative distribution function of two random variables X, Y is defined by  $F_{X,Y}(x,y) = \mathbb{P}(X \le x, Y \le y).$ X and Y are said to be jointly continuous with joint pdf  $f_{X,Y}$  if their cdf can be written

as an integral

With T(A) = B:

Joint distributions

While we can change  $f_{X,Y}$  at finitely many points without changing the integral, thus violating continuity, in general where  $F_{X,Y}$  is differentiable it is natural to write

 $F_{X,Y}(x,y) = \int_{-\infty}^{x} \int_{-\infty}^{y} f(u,v) du dv$ 

**Theorem 14** Suppose  $T:(x,y)\mapsto (u,v)$  is a bijection from some  $D\subseteq\mathbb{R}^2$  to some  $R \subseteq \mathbb{R}^2$ . We define the jacobian as

$$T(X,Y) \ are \ jointly \ continuous \ with \ joint \ pdf$$
 
$$f_{U,V}(u,v) = \begin{cases} f_{X,Y}(x(u,v),y(u,v))J(u,v) & if \ (u,v) \in R \\ 0 & otherwise \end{cases}$$

$$= \iint_A f_{X,Y}(x,y)\,\mathrm{d}x\,\mathrm{d}y$$
 
$$= \iint_B f_{X,Y}(x(u,v),y(u,v))J(u,v)\,\mathrm{d}u\,\mathrm{d}v.$$
 So the result is immediate via substitution.

 $\mathbb{P}((U,V) \in B) = \mathbb{P}((X,Y) \in A)$ 

Jacobian used is that of  $T^{-1}$  rather than that of T, so in fact the entire statement might be better expressed using  $T^{-1}$  than T.

and we can then define 
$$W_1, W_2, \ldots, W_n$$
 by 
$$\begin{pmatrix} W_1 \\ W_2 \\ \vdots \\ W_n \end{pmatrix} = A \begin{pmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_n \end{pmatrix} + \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix}.$$
 For  $A$  invertible then we can apply change of variables to get a joint distribution  $f_{\mathbf{W}_n}$  gi

For A invertible then we can apply change of variables to get a joint distribution  $f_{\mathbf{W}}$ , giving  $f_{\mathbf{W}}(\mathbf{w}) = \frac{1}{(2\pi)^{n/2}|\det A|} \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu})^{\top} (AA^{\top})^{-1}(\mathbf{w} - \boldsymbol{\mu})\right)$ 

## The result is immediate for transient chains, as with probability 1, $V_i(n)$ is bounded. If instead the chain is recurrent, take $R_k$ as the time between the kth and (k+1)th visits to i, which are i.i.d. with mean $m_i$ , and by the strong law of large numbers their sample mean tends to $m_i$ almost surely, indicating that where $T_k$ is the time of the kth

get  $\mathbb{E}V_n(i)/n \to 1/m_i$ .

visits to state i before time n, that is

Stationary distributions

distribution if

surely.  $V_i(n)/n$  is a bounded increasing sequence, so it is known to converge and by the previous statement it must converge to  $m_i$ . **Theorem 16** Let P be an irreducible transition matrix. Then P has a stationary

visit to i, as  $T_1$  is finite thus  $T_k/k \to m_i$  almost surely. We get that  $V_i(T_k) = k$ , so

 $V_i(T_k)/T_k = k/T_k \to 1/m_i$  almost surely as  $k \to \infty$ , and  $T_k \to \infty$  as  $k \to \infty$  almost

distribution if and only if P is positive recurrent, and the stationary distribution  $\pi$  is

same result. The converse is determined by taking the expected rate of visits  $\mathbb{E}V_n(i)/n$  for

 $X_0$  distributed by  $\pi$ , noting that this is  $\pi_i$ , and that by probabilistic convergence we can

**Theorem 17** If P is irreducible and aperiodic with stationary distribution  $\pi$ , then for

initially distributed by  $\pi$  with transition matrix also P. With  $T = \inf\{n \geq 0 \mid X_n = Y_n\}$ ,

we can consider  $W_n = (X_n, Y_n)$  as a markov chain, which is irreducible with a stationary

any initial distribution, for all 
$$i \in I$$
,  $\mathbb{P}(X_n = i) \to \pi_i$  as  $n \to \infty$ , and in particular for all  $i, j \in I$ ,  $p_{ij}^{(n)} \to \pi_j$ .

Let  $(X_n)$  be Markov distributed with initial distribution  $\lambda$ , transition matrix  $P$ , and  $(Y_n)$ 

distribution, so is positive recurrent and  $\mathbb{P}(T<\infty)=1$ . Thus we can define the chain  $Z_n$ as  $X_n$  for n < T and  $Y_n$  for  $n \ge T$ , and it turns out that this is Markov. Thus the result follows from here (**check**). Time reversal **Theorem 18** For P an irreducible transition matrix with stationary distribution  $\pi$ , and  $(X_0,\ldots,X_N) \sim \operatorname{Markov}(\pi,P)$ . Then for  $0 \leq n \leq N$ , with  $Y_n = X_{N-n}$ ,  $(Y_0,\ldots,Y_N) \sim \operatorname{Markov}(\pi,Q) \text{ with } Q = (q_{ij}) \text{ for } q_{ij}$ 

First we take the matrix Q, and observe that it is stochastic by taking the sum of each

 $= \mathbb{P}(X_N = i_0 | X_{N-1} = i_1, \dots, X_0 = i_N)$ 

 $\mathbb{P}(X_0 = i_N, \dots, X_{N-1} = i_1)$ 

 $= p_{i_1 i_0} \mathbb{P}(X_0 = i_N, \dots, X_{N-1} = i_1)$ 

 $\mathbb{P}(Y_0 = i_0, \dots, Y_N = i_N) = \mathbb{P}(X_0 = i_N, \dots, X_N = i_0)$ 

 $=\pi_{i_N}\prod_{k=1}^N p_{i_ki_{k-1}} \ =\pi_{i_N}\prod_{k=1}^N rac{\pi_{i_{k-1}}}{\pi_{i_k}}q_{i_{k-1}i_k}$ 

and Q also has stationary distribution  $\pi$ .

row.

and consequently we immediately get that  $\mathbb{P}(Y_0 = i) = \pi_i$ , as well as that  $\mathbb{P}(Y_n = j \mid Y_{n-1} = i_{n-1}, \dots, Y_0 = i_0) = q_{i_{n-1}j}$ , so independent of  $i_0, \dots, i_{n-2}$  and thus  $Y \sim \operatorname{Markov}(\pi, Q)$ .

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

**№T<sub>E</sub>X** TikZposter

 $f_{X,Y}(x,y) = \frac{\partial F_{X,Y}}{\partial x \partial y}(x,y)$ For suitably nice (Borel measurable) sets  $A \subseteq \mathbb{R}^2$ ,  $\mathbb{P}((X,Y) \in A) = \iint_A f_{X,Y}(x,y) \, \mathrm{d}x \, \mathrm{d}y$ We also get the obvious results of  $f_X(x) = \int_{\mathbb{R}} f_{X,Y}(x,y) \, dy$ ,  $f_Y(y) = \int_{\mathbb{R}} f_{X,Y}(x,y) \, dx$ .  $J = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial y}{\partial u} \\ \frac{\partial x}{\partial v} & \frac{\partial y}{\partial v} \end{vmatrix}$  $= \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial y}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix}$ If X, Y have joint pdf  $f_{X,Y}$  which is 0 outside D, then the random variables (U,V) =

I'm keeping the notation from lectures here, although in all honesty some weird choices were made here. For instance, the function 
$$(u, v) \mapsto (x(u, v), y(u, v))$$
 is just  $T^{-1}$ . The Jacobian used is that of  $T^{-1}$  rather than that of  $T$ , so in fact the entire statement might be better expressed using  $T^{-1}$  than  $T$ .

The above can then be generalised to the case of joint distributions of  $n > 2$  random

variables, for which the Jacobian becomes the determinant of an  $n \times n$  matrix. With

 $= rac{1}{(2\pi)^{n/2}} \exp\left(-rac{1}{2} oldsymbol{z}^{ op} oldsymbol{z}
ight)$ 

 $Z_1, Z_2, \ldots, Z_n$  standard normal variables, their joint density function can be written as

 $f_{\mathbf{Z}}(\mathbf{z}) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z_i^2}{2}\right)$ 

Let X be a markov chain with transition matrix P. A distribution over  $X_0, \pi$ , is a stationary

 $\pi P = \pi$ 

**Theorem 15 (Ergodic theorem)** Let P be irreducible. Let  $V_i(n)$  be the number of

so we have that if  $X_0$  is distributed by  $\pi$ , then so will  $X_n$  be for all  $n \geq 0$ .

$$V_i(n) = \sum_{r=0}^{n-1} \mathbb{1}(X_r = i).$$
 Then for any initial distribution, and for all  $i \in I$ , 
$$\frac{V_i(n)}{n} \to \frac{1}{m_i} \quad almost \; surely \; as \; n \to \infty$$

unique, given by  $\pi_i = 1/m_i$ . If P is positive recurrent we get that  $\pi_i = 1/m_i$  is an eigenvector immediately for finite state spaces, and for infinite we get an upper bound on  $\pi_i$  in terms of  $\pi_i$  which gives the

This follows immediately from the definitions. These equations are sometimes referred

**Theorem 19** Let P be an irreducible transition matrix with stationary distribution  $\pi$ . P is reversible iff for all  $i, j \in I$ 

to as the detailed balance equations.

We say that a transition matrix P is reversible if P = Q.

**Theorem 20** If the matrix P and the distribution  $\pi$  are in detailed balance, then  $\pi$ is stationary for P. If a drunk person was wandering with uniform random distribution around town, they This follows as  $\pi_j = \sum_i \pi_j p_{ji} = \sum_i \pi_i p_{ij}$  for any j. It is this characterisation of stationary distributions which makes time reversal so useful.

From this we get the notion of a mean return time,  $m_i = 1 + \sum_i p_{ij} k_i^{\{i\}}$ . If i is recurrent but  $m_i$  is infinite, we say that i is null recurrent. If however  $m_i < \infty$  then i is positive

using the same idea as for hitting probabilities.