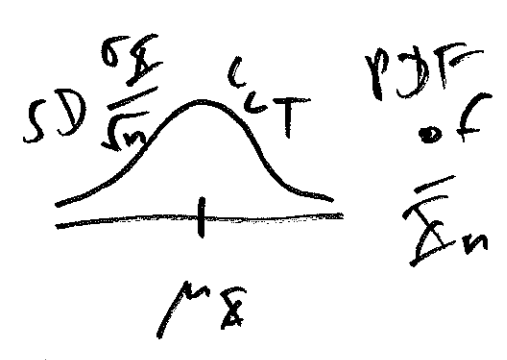the
Delta
Method

the CLT says that if $X_i \overset{IID}{\sim}$ (any) dist. with finite mean $\mu_X$ and finite variance $\sigma_X^2$, then

the distribution of $\dfrac{\overline{X}_n - \mu_X}{\sigma_X/\sqrt{n}}$ for large $n$

is approximately standard normal, where $\overline{X}_n = \dfrac{1}{n} \sum\limits_{i=1}^{n} X_i$.

---

this is equivalent to saying that



$$\overline{X}_n \sim N\left(\mu_X, \frac{\sigma_X^2}{n}\right)$$
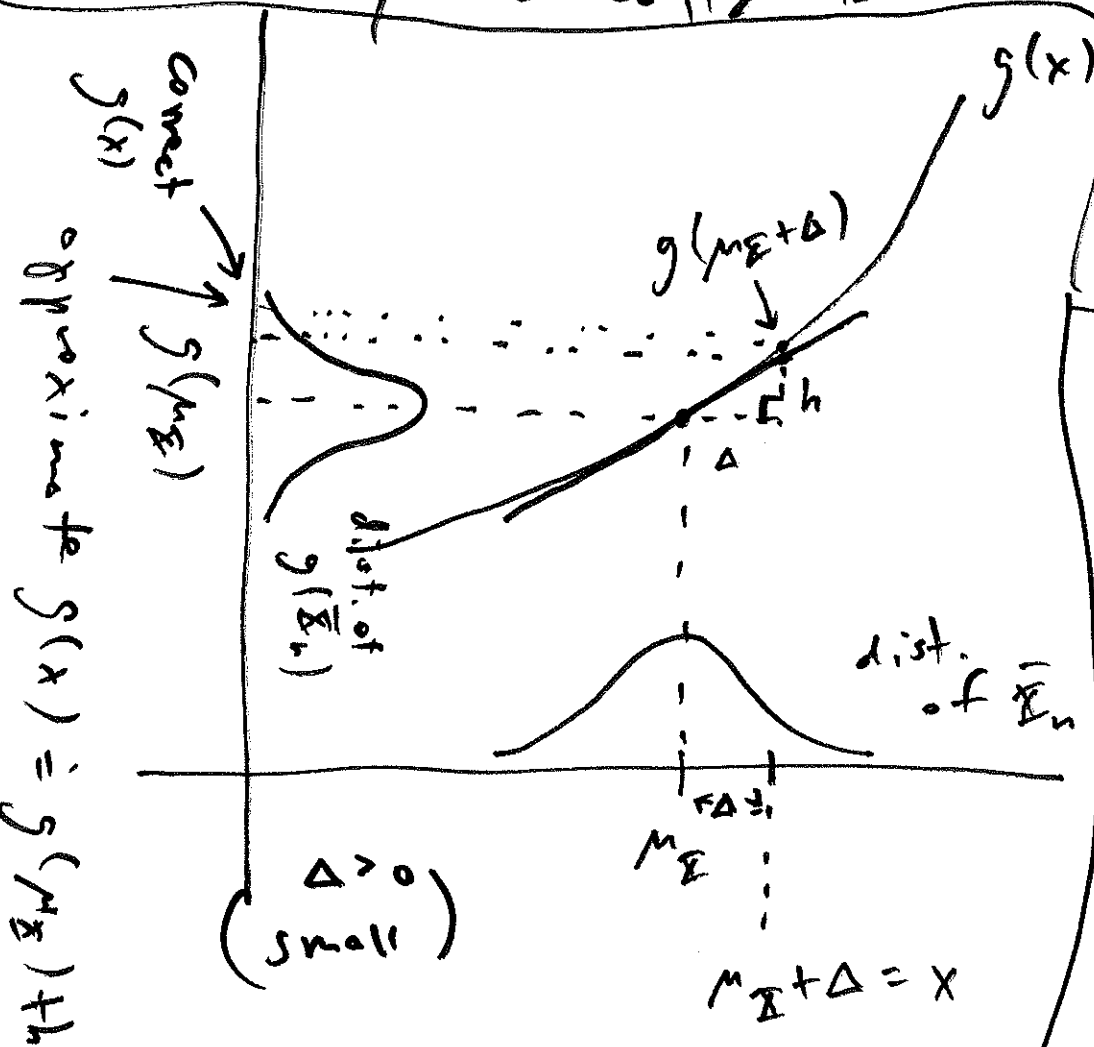
question: If $g(x)$ is a sufficiently "nice" function, is there a comparable result for $g(\overline{X}_n)$?

---

Answer: Yes, via a Taylor-series-based approach called the Delta method

$\bar{X}_n$ should be close to $\mu_{\bar{X}}$ for large $n$ (that's the (weak) Law of large Numbers); this suggests making a two-term Taylor expansion of $g(\bar{X}_n)$ around the point

$$x = \mu_{\bar{X}} : \quad g(\bar{X}_n) \doteq g(\mu_{\bar{X}}) + g'(\mu_{\bar{X}})(\bar{X}_n - \mu_{\bar{X}})$$

this is why it's called the $\Delta$ (Delta) - Method



$g(x)$

$g(\mu_{\bar{X}} + \Delta)$

correct $g(x)$

$g(\mu_{\bar{X}})$

dist. of $g(\bar{X}_n)$

approximate $g(x) \doteq g(\mu_{\bar{X}}) + g'(\mu_{\bar{X}})(x - \mu_{\bar{X}})$

dist. of $\bar{X}_n$

$\mu_{\bar{X}}$

$\Delta > 0$ (small)

$\mu_{\bar{X}} + \Delta = x$

so $\Delta = x - \mu_{\bar{X}}$

$$\frac{h}{\Delta} = g'(\mu_{\bar{X}})$$

so

$$g(x) \doteq g(\mu_{\bar{X}}) + h$$

$$= g(\mu_{\bar{X}}) + g'(\mu_{\bar{X}}) \cdot \Delta$$

$$= g(\mu_{\bar{X}}) + g'(\mu_{\bar{X}})(x - \mu_{\bar{X}})$$

$$g(\bar{X}_n) \doteq g(\mu_{\bar{X}}) + g'(\mu_{\bar{X}})(\bar{X}_n - \mu_{\bar{X}}) \quad \text{so}$$

constant ↑ r.v.

$$E[g(\bar{X}_n)] \doteq E[g(\mu_{\bar{X}}) + g'(\mu_{\bar{X}})(\bar{X}_n - \mu_{\bar{X}})]$$

$$= g(\mu_{\bar{X}}) + g'(\mu_{\bar{X}})\left[ E(\bar{X}_n)^{\nearrow 0} - \mu_{\bar{X}} \right]$$

$$\text{So} \quad E[g(\bar{X}_n)] \doteq g(\mu_{\bar{X}}) = g\left[ E(\bar{X}_n) \right]$$

and

constant

$$V[g(\bar{X}_n)] \doteq V\left[ g(\mu_{\bar{X}}) + g'(\mu_{\bar{X}})(\bar{X}_n - \mu_{\bar{X}}) \right]$$

r.v.

$$= [g'(\mu_{\bar{X}})]^2 \cdot V(\bar{X}_n - \mu_{\bar{X}})$$

$$\text{so} \quad V[g(\bar{X}_n)] \doteq [g'(\mu_{\bar{X}})]^2 V(\bar{X}_n)$$

$$V[g(\bar{X}_n)] \doteq [g'(\mu_{\bar{X}})]^2 \frac{\sigma_{\bar{X}}^2}{n}$$

There's one hidden assumption in

this calculation: $g'(\mu_{\bar{X}}) \neq 0$.

___

with finite variance

This works for any r.v., not just $\bar{X}_n$:

___

$V$ any r.v. with finite variance $\sigma_V^2$ (and

therefore finite mean $\mu_V$), $W = g(V)$

$\rightarrow E(W) \doteq g(\mu_V)$ and

$$V(W) \doteq [g'(\mu_V)]^2 \sigma_V^2,$$

$\boxed{\Delta \text{ method}}$ $\underset{\text{1}}{\textcircled{\text{part}}}$

provided $g'(v)$ is continuous and

$g'(\mu_V) \neq 0$

Moreover, if $V \sim$ Normal

then $W = g(V) \sim$ Normal

also

$\boxed{\Delta \text{ method}}$ $\textcircled{\text{part 2}}$

**Example** A bank typically has a single queue (line) at which customers arrive to transact banking business.

Let $X_i$ = time customer $i$ waits from reaching the head of the queue until served. To be completely realistic, the dist. of $X_i$ would vary by day of week and time of day, so pick a single time slot (e.g. Tue 10-10.15am) and observe the $X_i$ from week to week only in that time slot; now the $\{X_i, i = 1, 2, ...\}$ form a _stationary_ stochastic process with fixed (non-time-varying) finite $E(X_i) = \mu_X$

$> 0$

and fixed (non-time-varying) finite

$V(\underline{X}_i) = \sigma_X^2$.

Gather data over many weeks and form $\bar{X}_n = \frac{1}{n}\sum_{i=1}^{n} X_i$ for large $n$.

The rate of service is defined to be $g(\mu_X) = \frac{1}{\mu_X}$, which would naturally be estimated by $g(\bar{X}_n) = \frac{1}{\bar{X}_n}$.

complication: seasonal effects (ignored here)

| $E(\bar{X}_n) = \mu_X$ $V(\bar{X}_n) = \frac{\sigma_X^2}{n}$ | $g(x) = \frac{1}{x} = x^{-1}$ | $g'(x) = -\frac{1}{x^2}$ |
|---|---|---|
| | $g'(\mu_X) = -\frac{1}{\mu_X^2}$ | $\bar{X}_n \sim$ Normal by CLT |

so $\Delta$-method says $g(\bar{X}_n) = \frac{1}{\bar{X}_n} \sim$ Normal also,

with mean $g(\mu_X) = \frac{1}{\mu_X}$ and variance $\sigma_X^2/(n\mu_X^4)$

$\left(g'(\mu_X)\right)^2 = \frac{1}{\mu_X^4} \neq 0$

| Specific calculation | Under some plausible assumptions, we've seen that $(X_i | \lambda) \overset{IID}{\sim} \text{Exponential}(\lambda)$ |
|---|---|

may be a reasonable model for waiting times.

$E(X_i) = \dfrac{1}{\lambda}$, $V(X_i) = \dfrac{1}{\lambda^2}$  $(X_i | \lambda)$ has PDF

$\quad = \mu_X$, $\quad = \sigma_X^2$

$$f_{X_i}(x_i | \lambda) = \lambda e^{-\lambda x_i} I(x_i > 0)$$

So $\dfrac{1}{\bar{X}_n}$ should (large $n$) be approximately Normal with mean $\dfrac{1}{\frac{1}{\lambda}} = \lambda$

and SD $\dfrac{\sigma_X}{\mu_X^2 \sqrt{n}} = \dfrac{\frac{1}{\lambda}}{\left(\frac{1}{\lambda}\right)^2 \sqrt{n}} = \dfrac{\lambda}{\sqrt{n}}$.

(discrete or continuous)

| Fancy version of $\Delta$-method | $q_1, q_2, \ldots$ sequence of r.v.; $F^*$ continuous cdf; |
|---|---|

$\theta$ a real number; $a_1, a_2, \ldots \uparrow \infty$ positive sequence

$g(\cdot)$ a ^(real-valued) function of a real variable  (329)

such that $g'(\cdot)$ is continuous and

$g'(\theta) \neq 0$; then if $a_n(\mathcal{I}_n - \theta) \xrightarrow{D} F^*$,

$$a_n\left[\frac{g(\mathcal{I}_n) - g(\theta)}{|g'(\theta)|}\right] \xrightarrow{D} F^* \quad \text{also}$$

Typical application:

$\mathcal{I}_1, \mathcal{I}_2, \ldots$ IID

$$\mathcal{I}_n = \bar{X}_n = \frac{1}{n}\sum_{i=1}^{n} X_i \; ; \; \theta = \mu_X \; ; \; a_n = \frac{\sqrt{n}}{\sigma_X} \; ;$$

$F^* = \Phi$, the standard normal CDF.

In this context the theorem says that

if $\dfrac{\bar{X}_n - \mu_X}{\sigma_X/\sqrt{n}} \sim N(0,1)$ $\xrightarrow{\text{then}}$ $\dfrac{g(\bar{X}_n) - g(\mu_X)}{|g'(\mu_X)|\sigma_X/\sqrt{n}}$

~~(28 Aug 17)~~ is also $\sim N(0,1)$

A little bit more about the continuity correction | Tay-Sachs case study, revisited

$X = \#$ T-S babies in family of $n=5$ children, both parents carriers so that

$$P(\text{T-S baby}) = \frac{1}{4} = p \quad \boxed{X \sim \text{Binomial}(n,p)}$$

But also let $T_i = \begin{cases} 1 & \text{if child } i \text{ is T-S baby} \\ 0 & \text{else} \end{cases}$

$i = 1, \ldots, n = 5$

Then $(T_i) \stackrel{IID}{\sim} \text{Bernoulli}(p)$ and $\underline{X} = \sum_{i=1}^{n} T_i$

$(i = 1, \ldots, n)$

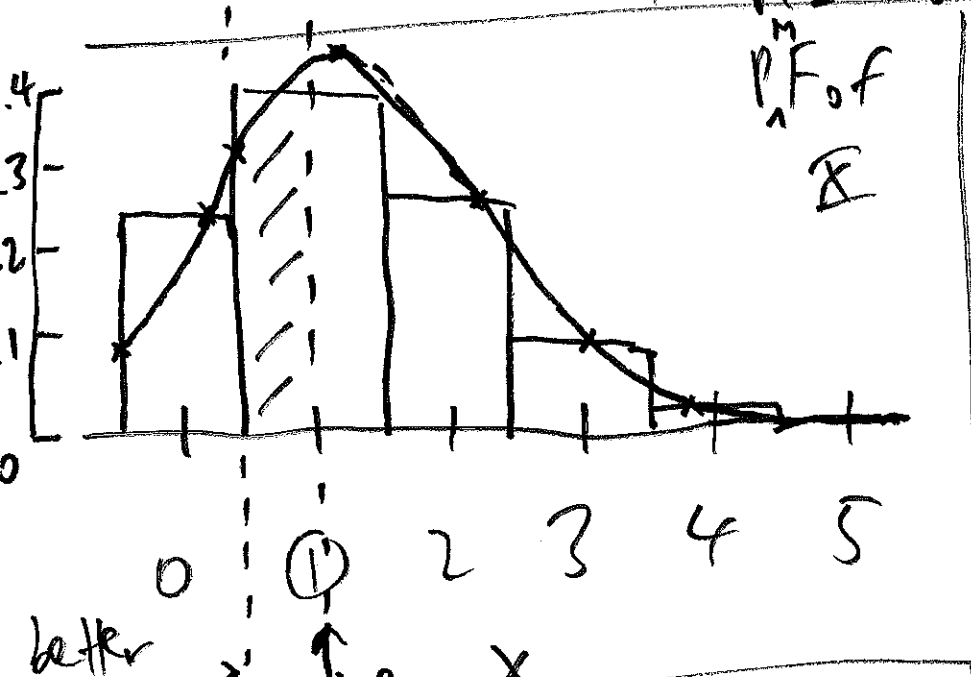So by the CLT the dist. of $X$ should be approximately Normal with mean

$\mu_X = E(X) = np = 1.25$ and SD

$$\sigma_{\underline{X}} = \sqrt{V(\underline{X})} = \sqrt{np(1-p)} \doteq 0.98$$

on day 1 of this class we worked out

that $p(1 \text{ or more } T\text{-}s \text{ ladies}) = P(\underline{X} \geq 1)$

$1 - p(\text{no } T\text{-}s \text{ ladies}) = 1 - (1-p)^n \doteq 0.76$
   $= 1 - P(\underline{X}=0)$



PMF of $\underline{X}$

better approx → naive approx   $X$

$P(\underline{X} \geq 1) \doteq 1 - P(\underline{X}' < 1)$

$= 1 - 0.398$

$= 0.602$ (quite a bad approximation)

Naive Normal approximation, from CLT:

SD 0.98   PDF of $\underline{X}'$

.398   .602

1.0
1.25

$$\frac{1.0 - 1.25}{0.98} \doteq -0.26$$

Improved approximation obtained by
paying attention to the edges of the
histogram ( $\overset{M}{\underset{\wedge}{PF}}$ ) bars:

---

Normal approximation

with continuity correction

$$P(X \geq 1) \doteq 1 - P(X' < 0.5)$$

$$\doteq 1 - .219$$

$$\doteq 0.781 \quad \left( \begin{array}{l} \text{correct answer } 0.76; \\ \text{much better approx.} \end{array} \right)$$

SD 0.98          PDF

.2192      .781        of
                            $X'$

0.5  1.25

$$\frac{0.5 - 1.25}{0.98} \doteq -0.77$$

---

| Markov Chains | Recall the definition of a stochastic process : |

Def.) A sequence of rvs $X_1, X_2, \ldots$
is called a <u>stochastic process</u> with
<u>discrete time parameter</u> $t = 1, 2, \ldots$.

$X_1$ is the <u>initial state</u> of the process;
$X_n$, $n \geq 1$ is the <u>state</u> of the process
at time $t = n$. | The simplest possible
discrete-time stochastic process is
an <u>IID</u> sequence of rvs $(X_1, X_2, \ldots)$.

Suppose that there's a <u>parameter</u> $\theta$
such that $(X_i \mid \theta) \overset{IID}{\sim}$ from some dist.
depending on $\theta$. | <u>Q:</u> Does this process
have a memory?

Example, revisited | Machine with $\theta$ dial from 0 to 1, produces IID Bernoulli($\theta$) trials $\underline{X}_i$. | The process $(\underline{X}_1, \underline{X}_2, ...)$ does have a memory, for you, if $\theta$ is unknown to you: the information that 17 out of the first 20 trials were successes helps you to predict $\underline{X}_{21}$, because it's reasonable to conclude from $\underline{X}_1, ..., \underline{X}_{20}$ that $\theta$ is around $\frac{17}{20} \doteq 0.85$, so $\underline{X}_{21}$ will be probably a success. | But the process

$$\{(\underline{X}_i | \theta), i = 1, 2, ..\}$$ has no memory once $\theta$ is known: information about

the first $n$ trials is irrelevant to
your prediction of $X_{n+1}$ if you know

$\theta$. $\sqrt{\phantom{xx}}$ An IID process $(X_i | \theta) \overset{IID}{\sim}$

is called a white-noise (stochastic)

process or a white noise time series.

Q: what's the next level of complexity
for discrete-time stochastic processes
up from white noise? A: Allow $X_{n+1}$

to depend on $X_n$ but not on $X_{n-1}, X_{n-2}, \ldots$

(i.e., let the process have a short-term
memory, ① time period back in the
past).

From now on, I'll suppress the dependence of the process on $\theta$ in the notation.

**Def.** A $\underset{\wedge}{\text{discrete-time}}$ stochastic process is a (first-order) <u>Markov chain</u> if for $n = 1, 2, \ldots$; $b$ any real number; and for all possible sequences of states $x_1, x_2, \ldots$

$$P(X_{n+1} \leq b \mid X_1 = x_1, \ldots, X_n = x_n)$$
$$= P(X_{n+1} \leq b \mid X_n = x_n).$$

In other words, the only thing you need to know to simulate where the Markov chain is going <u>next</u> is <u>where it is now</u>.

(Can define higher-order Markov chains
with memory of 2 or more time periods;
we won't pursue that here.)

| Def.

The set of values ~~the~~ a Markov chain can

take on is called its _state space_

$S$, which may be finite or infinite.

( Can also have Markov chains unfolding
in continuous time, e.g. $X_t =$ stock
price at time $t =$ seconds, milliseconds,
microseconds, ...; we also won't pursue
that here.)

| It's easy to write down
the joint PMF of a Markov chain with finite
$S$:

Consequences

**Def.** A Markov chain with a finite state space is called a <u>finite</u> Markov chain.

**Def.**

① $(X_1, X_2, ...)$ finite Markov chain →

$$P(X_1 = x_1, ..., X_n = x_n) =$$

$$P(X_1 = x_1) \cdot P(X_2 = x_2 \mid X_1 = x_1) \cdot$$

$$P(X_3 = x_3 \mid X_2 = x_2) \cdot \cdots$$

$$P(X_n = x_n \mid X_{n-1} = x_{n-1}).$$

Suppose you have a finite Markov chain with $k$ possible states numbered $1, ..., k$ ($k$ integer $\geq 2$) → $\{ P(X_{n+1} = j \mid X_n = i),$ $i, j = 1, ..., k, n = 1, 2, ... \}$ are called the <u>transition</u> distribution of the Markov chain.

If $p(\underline{X}_{n+1}=j \mid \underline{X}_n=i)$ is the same for all $n$, the transition distribution is said to be **stationary** (D5). (time-homogeneous)

If the Markov chain does have a stationary transition distribution, then the probabilities

$$p_{ij} \triangleq p(\underline{X}_{n+1}=j \mid \underline{X}_n=i)$$ completely characterize the Markov chain's behavior.

Can arrange the $p_{ij}$ in a matrix called the transition matrix.

to state $p_{ij}$

$$p = \sum_{k=k}^{k} \quad \text{(from)} \quad \text{state } k \begin{array}{c} 1 \\ 2 \\ \vdots \\ k \end{array}\begin{bmatrix} p_{11} & p_{12} & \cdots & p_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & & \ddots & \vdots \\ p_{k1} & p_{k2} & \cdots & p_{kk} \end{bmatrix}$$

All of the elements of $\underline{\underline{P}}$ are non-negative (they're probabilities), and all of the row sums are 1 (because the chain has to go _somewhere_), i.e.

$$\sum_{j=1}^{k} p_{ij} = 1 \text{ for all } i = 1, \ldots, k.$$

Def.

_____
_matrix versus quaternion_

A square matrix $\underset{k = k}{P}$ with non-negative entries and $\overset{all}{\underset{\wedge}{\text{row}}}$ sums equal to 1 is called a _stochastic matrix_.

~~(Doubly stochastic)~~

Example $\Big\{$ $\underset{\wedge}{\text{Gene}}$ inheritance is Markovian: all we need to know to predict you is the genetic story of your parents

(your grand parents, ..., are irrelevant).

Suppose that

A gene of interest to you has two

alleles, A and a | Then a _state_ in

the Markov chain is of the form

$$\{ \substack{\text{allele 1} \\ \text{from} \\ \text{parent} \\ 1} \quad \substack{\text{allele 2} \\ \text{from} \\ \text{parent} \\ 1}, \quad \substack{\text{allele 1} \\ \text{from} \\ \text{parent} \\ 2} \quad \substack{\text{allele} \\ 2 \\ \text{from} \\ \text{parent 2}} \}, \quad \text{for}$$
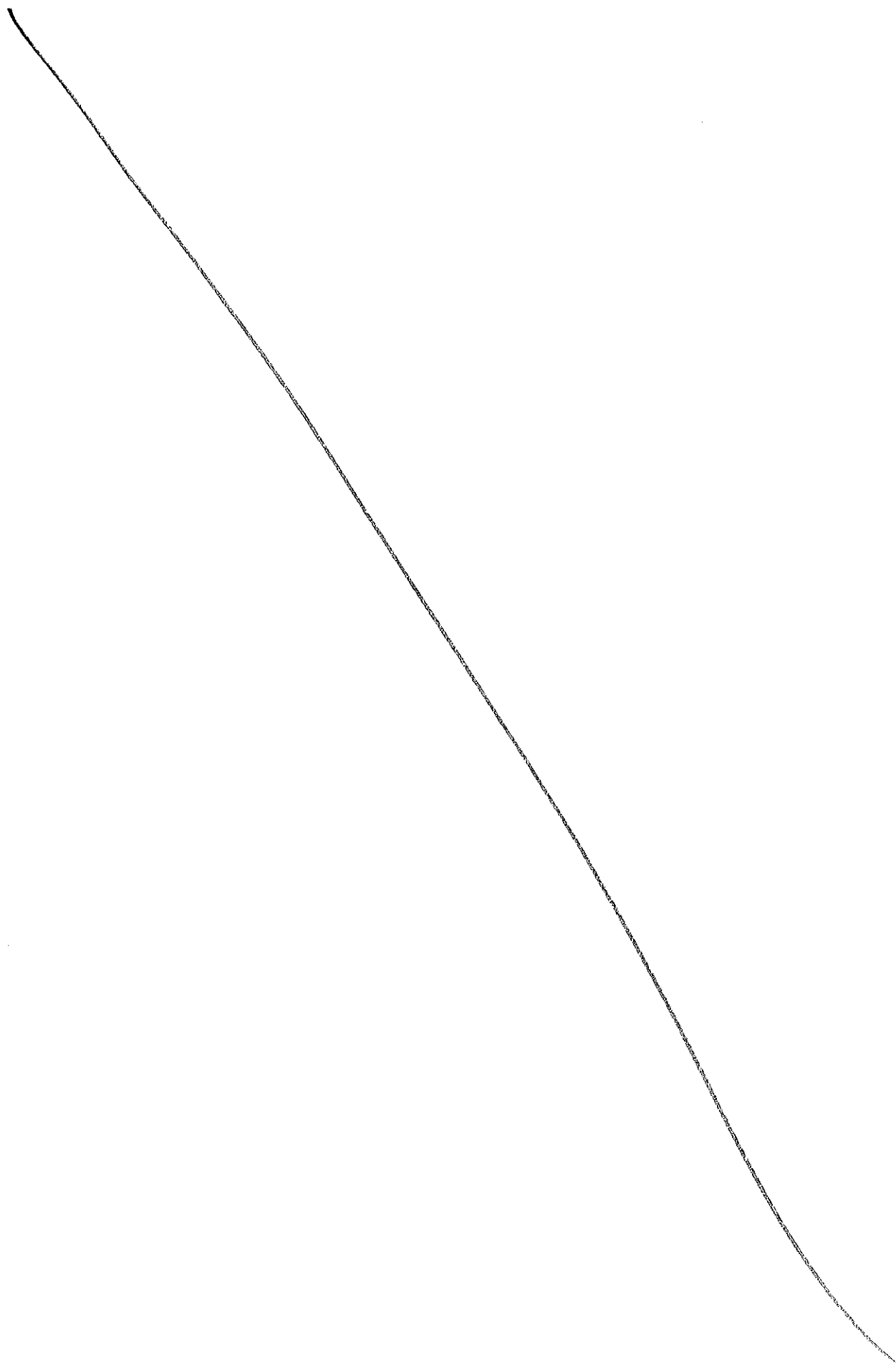
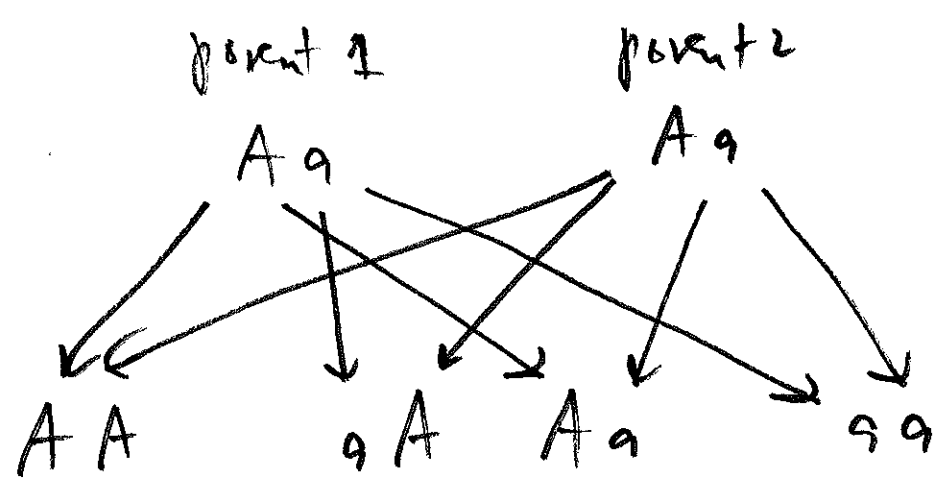example $\{Aa, Aa\}$. | Ignoring order

(because it's irrelevant in inheritance),

there are 6 possible states: $\{AA, AA\}$

$\{AA, Aa\}$, $\{AA, aa\}$, $\{Aa, Aa\}$, $\{Aa, aa\}$

and $\{aa, aa\}$.

parent 1 — Aa   parent 2 — Aa

AA   aA   Aa   aa

One possible inheritance sequence

---

offspring gets A or a from parent 1 and A or a (independently) from parent 2, (A or a) each with probability $\frac{1}{2}$.

Transition matrix

| from ↓ \ to → | {AA,AA} | {AA,Aa} | {AA,aa} | {Aa,Aa} | {Aa,aa} | {aa,aa} |
|---|---|---|---|---|---|---|
| {AA,AA} | 1 | 0 | 0 | 0 | 0 | 0 |
| {AA,Aa} | $\frac{1}{4}$ | $\frac{1}{2}$ | 0 | $\frac{1}{4}$ | 0 | 0 |
| {AA,aa} | 0 | 0 | 0 | 1 | 0 | 0 |
| {Aa,Aa} | $\frac{1}{16}$ | $\frac{1}{4}$ | $\frac{1}{8}$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{16}$ |
| {Aa,aa} | 0 | $\frac{1}{8}$ | 0 | $\frac{1}{4}$ | $\frac{1}{2}$ | $\frac{1}{4}$ |
| {aa,aa} | 0 | 0 | 0 | 0 | 0 | 1 |

Example (random walk) You're watching

a particle move around on the

integers $S = \{ \dots, -2, -1, 0, 1, 2, \dots \}$

over time: here are the rules:

wherever it is at time $t = n$,

it moves left 1 unit with prob $p_1$,

—— right 1 unit —————— $p_3$,

and it stays where it is with prob $p_2$,

where $0 < p_i < 1$ and $\sum_{i=1}^{3} p_i = 1$ | This is

clearly a Markov chain (why?);

what is its transition matrix?

$$
\underline{\underline{P}} =
\begin{array}{c|ccccc|c}
\text{from} \backslash \text{to} \to & \cdots & -2 & -1 & 0 & 1 & 2 & \cdots \\
\hline
\vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
-2 & \cdots & p_2 & p_3 & 0 & 0 & 0 & \cdots \\
-1 & \cdots & p_1 & p_2 & p_3 & 0 & 0 & \cdots \\
0 & \cdots & 0 & p_1 & p_2 & p_3 & 0 & \cdots \\
1 & \cdots & 0 & 0 & p_1 & p_2 & p_3 & \cdots \\
2 & \cdots & 0 & 0 & 0 & p_1 & p_2 & \cdots \\
\vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots
\end{array}
$$

This is an example of a band matrix, in which the only non-zero entries are on the main diagonal and 1 diagonal either way from the main diagonal; since there are only 3 non-zero diagonals, $\underline{\underline{P}}$ is said to be tridiagonal.

Moreover, all of the main diagonal entries are the same ($p_2$); all of the entries 1 diagonal below are also the same ($p_1$); and all of the entries 1 diagonal above are also the same ($p_3$).

Such matrices are called toeplitz (named after Otto Toeplitz, (1881-1940) a German mathematician who was fired by the Nazis from his university position in 1935 for being Jewish!) (died of tuberculosis at 58)

Q: Start this process, which is called a random walk, at 0 & let it go; where is the particle likely to be at time $n$, $n$ large?

A:| Suppose, for example, that $(p_1, p_2, p_3) = (0.1, 0.3, 0.6)$. Then you would expect the particle to drift off to $+\infty$. Similarly, $(p_1, p_2, p_3) = (0.5, 0.25, 0.25)$ should yield a drift to $-\infty$. $(p_1, p_2, p_3) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$?

Can show that as $n \to \infty$ every integer is visited infinitely many times, and the expected time you must wait for the chain to return to $0$ (having started there) is also infinite.

The infinite random walk evidently has "too much freedom" to move around to get interesting results; let's bound it.