

ex. $E(X) = 1$, X non-negative \rightarrow (30)

$$P(X \geq 100) \leq \frac{1}{100}$$

The inequality is

sharp, meaning that the upper bound

$\frac{E(X)}{t}$ on $P(X \geq t)$ is attainable, \otimes

ex. $E(X) = 1$, X - nonnegative \rightarrow
put probability 0.99 on $X = 0$ and
0.01 on $X = 100$

\otimes but most of the time (i.e., for most distributions) it's a crude upper bound.
(30 May 19)

Can apply Markov inequality to the
rv. $Y = (X - E(X))^2$ to get

Chebyshev Inequality } X r.v. with $V(X)$ existing (302)
→ for every $t \geq 0$,

$$P\left[|X - E(X)| \geq t\right] \leq \frac{V(X)}{t^2}$$

(attributed to

Pafnuty Chebyshev (1821 - 1894), also a Russian mathematician, one of whose Ph.D. students was Markov)

Ex.

$$E(X) = \mu$$
$$V(X) = \sigma^2$$

Chebyshev says $P\left[\left|\frac{X - \mu}{\sigma}\right| \geq 3\right] \leq \frac{1}{3^2} = \frac{1}{9}$,

So no more than $\frac{1}{9} = 11\%$ of the probability in any distribution, with finite variance, can

be more than 3 SDs away from the mean (recall for Normal dist. this prob. is 0.3%)

This upper bound is also sharp, but 3.3
for most distributions it's (also) crude
(as with the Markov bound). Back to \bar{X}_n

$X_i \stackrel{i.i.d.}{\sim}$ some dist. with mean $E(X_i) = \mu$
($i=1, \dots, n$) and variance $V(X_i) = \sigma^2 < \infty$

~~We~~ ^{have} already shown that if $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

then $E(\bar{X}_n) = \mu$ for all $n=1, 2, \dots$

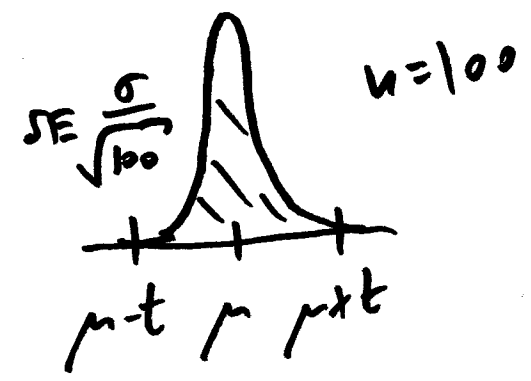
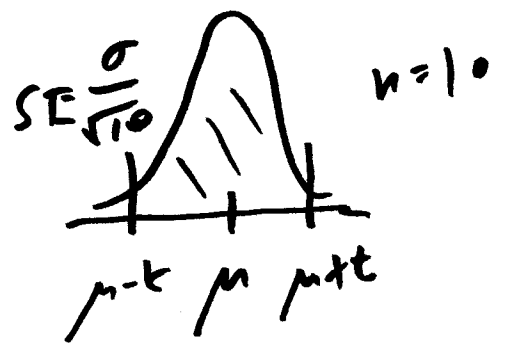
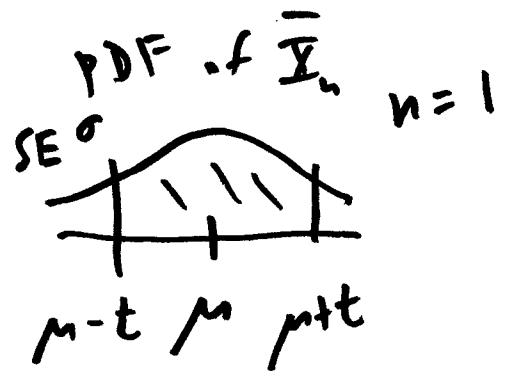
and $V(\bar{X}_n) = \frac{\sigma^2}{n}$.

Chebyshev then

$$\text{gives } P(|\bar{X}_n - \mu| \geq t) \leq \frac{\sigma^2}{nt^2} \text{ for all } t > 0$$

this can be

rewritten $P(|\bar{X}_n - \mu| < t) \stackrel{\epsilon > 0}{\geq} 1 - \frac{\sigma^2}{nt^2}$
 $\leftarrow \epsilon$



⋮

This suggests a way 304
 to quantify how close
 a r.v. like \bar{X}_n is to
 a constant like μ :

Def. A sequence Z_1, Z_2, \dots
 of r.v. is said to
 converge in probability
 to a constant b if

for all $\epsilon > 0$, $\lim_{n \rightarrow \infty} P(|Z_n - b| < \epsilon) = 1;$

this is denoted $Z_n \xrightarrow{P} b$

\nearrow
 An immediate
 consequence of this
 definition is

consequence of Chebyshev & this definition is

(weak)
Law of
Large
Numbers

$X_i \stackrel{\text{IID}}{\sim}$ a dist. with mean μ and variance $\sigma^2 < \infty$, $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$
(\bar{X}_n is consistent for μ)

$$\bar{X}_n \xrightarrow{P} \mu$$

This result has

the Italian mathematician

a long history: Gerolamo Cardano (1501-1576)

asserted it without proof; Jacob Bernoulli (1659-1705)

proved it for $X_i \sim \text{Bernoulli}(\theta)$

(it took him 20 years to find a correct

proof, published posthumously in 1713;

Bernoulli thought that this theorem proved

the existence of God); Siméon Denis Poisson

named it the Law of Large Numbers in

1837.

Corollary

If $Z_n \xrightarrow{P} b$ and $g(z)$

is continuous at $z=b$ then $g(Z_n) \xrightarrow{P} g(b)$.

Central Limit Theorem (CLT)

Example

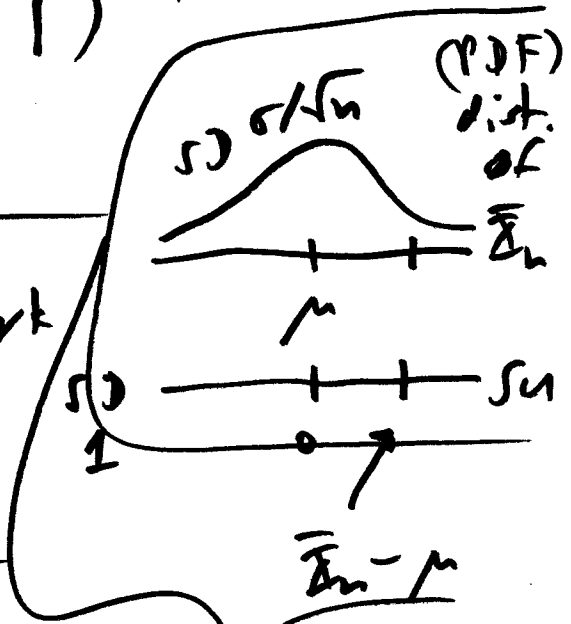
$X_i \sim \text{IID } N(\mu, \sigma^2)$, $\sigma < \infty$
($i=1, \dots, n$)

we know

that $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ has mean μ ,

variance $\frac{\sigma^2}{n}$ and is normally distributed,

so that $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$ for all $n=1, 2, \dots$



Q:

Does something like this work for other choices of

$X_i \sim \text{IID } [?]$

A: Yes: it's the most famous result in all of probability.

Central Limit Theorem (CLT)

$X_i \sim \text{IID}$ (any) dist. with mean μ and finite variance $0 < \sigma^2 < \infty$,

$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \rightarrow$

for large n : $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$

Careful statement) Def. X_1, X_2, \dots a sequence (30)

of r.v.; let F_n be the CDF of X_n

+ if there exists a CDF F^* such
that $\lim_{n \rightarrow \infty} F_n(x) = F^*(x)$ for all x at

which $F^*(x)$ is continuous, then

people say that $X_n \xrightarrow{D} F^*$ (X_n converges in distribution to F^*)

CLT) $X_j \stackrel{i.i.d.}{\sim}$ (any) dist. with mean μ
and variance $0 < \sigma^2 < \infty$, $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

$$\rightarrow \sqrt{n} \left(\frac{\bar{X}_n - \mu}{\sigma} \right) \xrightarrow{D} N(0, 1)$$

the
CLT

also has a long history: it was

first demonstrated for $X_i \sim \text{IID Bernoulli}(p)$
by the French/British mathematician
Abraham de Moivre (1667 - 1754) in
1733; almost forgotten until revived by
the French mathematician Pierre-Simon de
Laplace (1749 - 1827) in 1812; almost
forgotten again until 1901, when the
Russian mathematician Aleksandr Lyapunov
gave a more general proof; ^{even} more general
proof provided by JW Lindeberg (Finnish
mathematician (1876 - 1932)) and independently
by Paul Lévy (French mathematician (1886 -
1971)) in the early 1920s.

(26 Aug 19)
CLT name due to
(1882-1985) George Pólya in
Hungarian-American mathematician 1920

Example Contaminated water supply: (309)

X = arsenic concentration

Y = lead concentration
(same units) (both 30)

Interest focuses

$$R = \frac{Y}{X+Y}$$

(proportion of contamination due to lead)

$E(R) = E\left(\frac{Y}{X+Y}\right)$ difficult to calculate.

(IID)
Simulation approach Randomly sample n pairs (X_i, Y_i) from the joint PDF of (X, Y) , calculate $R_i = \frac{Y_i}{X_i + Y_i}$ and

$$\bar{R}_n = \frac{1}{n} \sum_{i=1}^n R_i \leftarrow \text{good Monte Carlo}$$

(Simulation) estimate of $E(R)$.

Q: How big does n need to be to achieve a desired accuracy target? (310)

By definition

$$|R_i| = \left| \frac{I_i}{\bar{X}_i + I_i} \right| \leq 1; \text{ can show that}$$

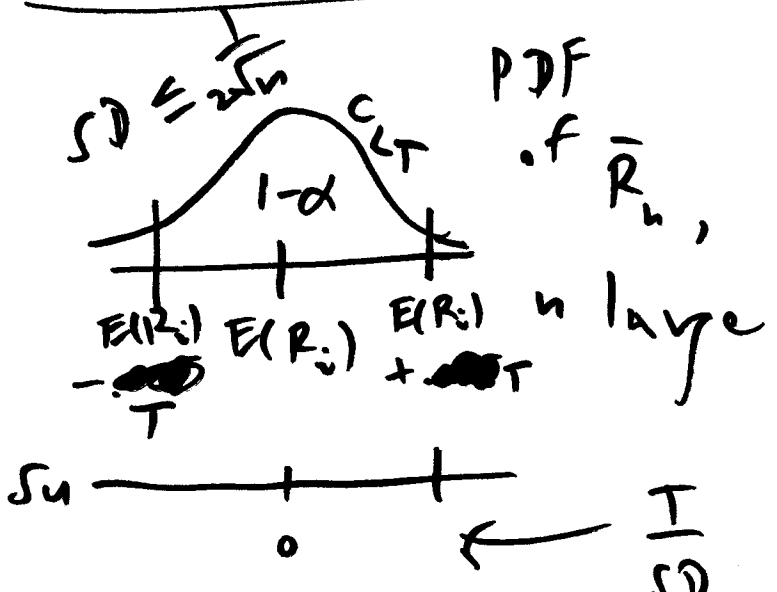
as a result $V(R_i) \leq \frac{1}{4}$.

CLT

Says that dist. of \bar{R}_n will be close to Normal for large n , with mean $E(R_i)$

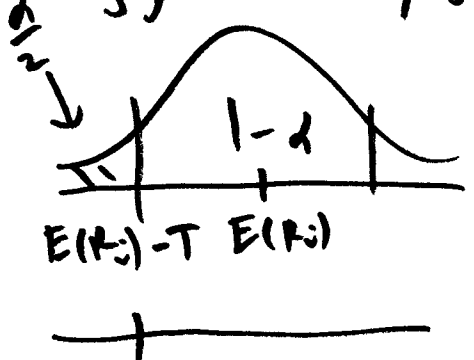
and variance $\frac{V(R_i)}{n} \leq \frac{1}{4n}$

Suppose we want \bar{R}_n to



differ from $E(R_i)$ by no more than one tolerance T with probability at least $(1-\alpha) \dots$

$SD \leq \frac{1}{2\sqrt{n}}$, so $\frac{1}{SD} \geq 2\sqrt{n}$ and



$\frac{-T}{SD} \leq 2T\sqrt{n}$

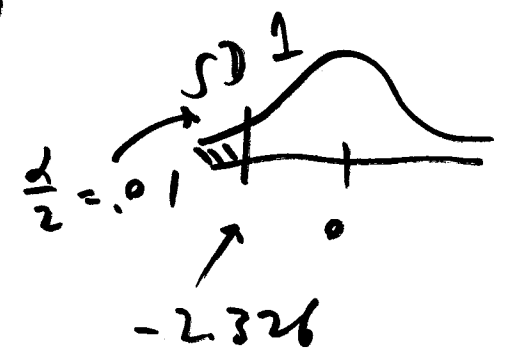
$$\Phi^{-1}\left(\frac{\alpha}{2}\right) = \frac{[E(R_i) - T] - E(R_i)}{SD} = \frac{-T}{SD} \leq 2T\sqrt{n}$$

from which $n \geq \left[\frac{\Phi^{-1}\left(\frac{\alpha}{2}\right)}{2T} \right]^2$

For instance, set $T = 0.005$ ($\frac{1}{2}$ of 1%)

and $d = .02$ to get

$$n \geq \left[\frac{-2.326}{2(.005)} \right]^2 = 54,119$$



simulation replications needed

Case Study: Escalators

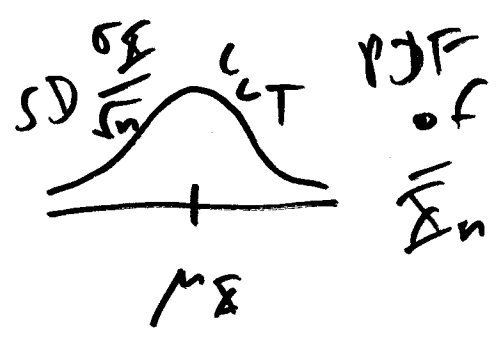
in the London Underground (👤)

the Delta Method

The CLT says that if $X_i \stackrel{i.i.d.}{\sim}$ (any) dist. with finite mean μ_X and finite variance σ_X^2 , then

The distribution of $\frac{\bar{X}_n - \mu_X}{\sigma_X/\sqrt{n}}$ for large n is approximately normal, where $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$.

This is equivalent to saying that



$$\bar{X}_n \sim N\left(\mu_X, \frac{\sigma_X^2}{n}\right)$$

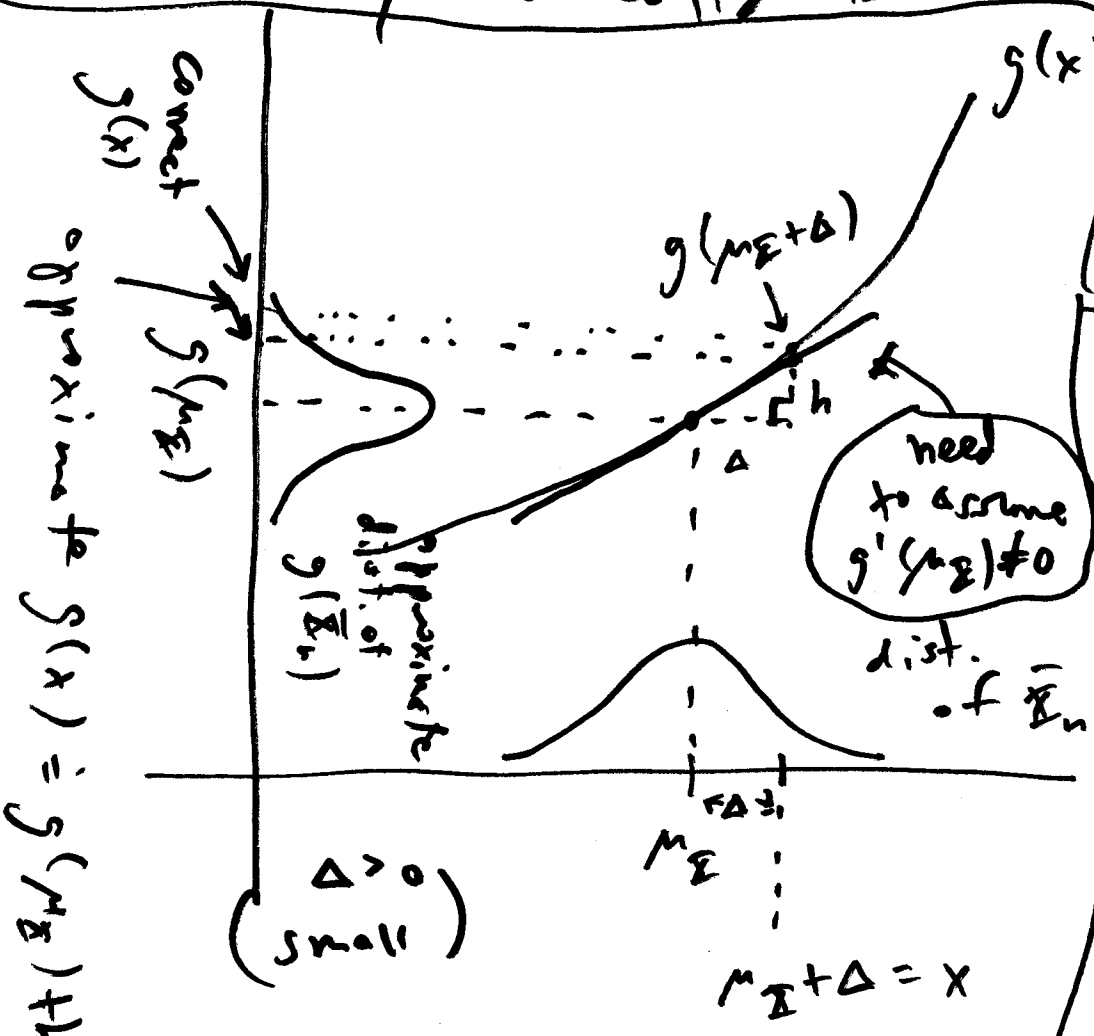
Question: If $g(x)$ is a sufficiently "nice" function, is there a comparable result for $g(\bar{X}_n)$?

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

\bar{x}_n should be close to μ_{Σ} 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

$$x = \mu_{\Sigma} : g(\bar{x}_n) \approx g(\mu_{\Sigma}) + g'(\mu_{\Sigma})(\bar{x}_n - \mu_{\Sigma})$$

this is why it's called the Δ (Delta) - method



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

so

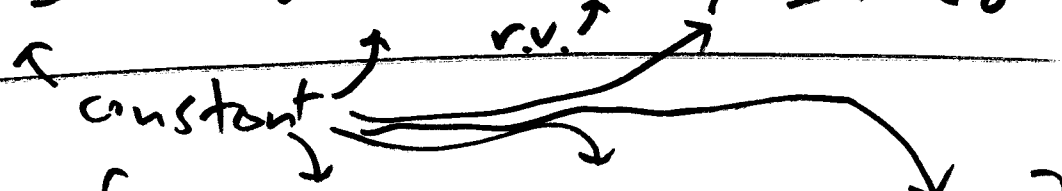
$$g(x) \approx g(\mu_{\Sigma}) + h$$

$$= g(\mu_{\Sigma}) + g'(\mu_{\Sigma}) \Delta$$

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

so $\Delta = x - \mu_{\Sigma}$

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



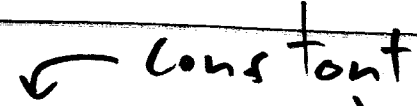
$$E[g(\bar{X}_n)] = E[g(\mu_X) + g'(\mu_X)(\bar{X}_n - \mu_X)]$$

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

$$\text{so } E[g(\bar{X}_n)] = g(\mu_X) = g[E(\bar{X}_n)]$$

and

$$V[g(\bar{X}_n)] = V[g(\mu_X) + g'(\mu_X)(\bar{X}_n - \mu_X)]$$



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



$$\text{so } V[g(\bar{X}_n)] = [g'(\mu_X)]^2 V(\bar{X}_n)$$

$$\text{i.e., } V[g(\bar{X}_n)] = [g'(\mu_X)]^2 \frac{\sigma_X^2}{n}$$

There's one hidden assumption in this calculation: $g'(\mu_X) \neq 0$.

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

Any $r.v.$ with finite variance σ_V^2 (and therefore finite mean μ_V), $W = g(V)$

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

$V(\bar{W}) = [g'(\mu_V)]^2 \sigma_V^2$, Δ method
part 1

provided $g'(v)$ is continuous and

$g'(\mu_V) \neq 0$

Moreover, if V is Normal then $W = g(V)$ is Normal also

Δ method part 2

Example | A bank typically has a 316
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 $\{X_i, i=1, 2, \dots\}$
form a stationary stochastic process
with fixed (non-time-varying) ^{finite} $E(X_i) = \mu_X$

70

and fixed (non-time-varying) finite (317)

$$V(\underline{X}_i) = \frac{\sigma^2}{n}$$

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

Complication:
seasonal effects
(ignored here)

is defined to be $g(\mu_X) = \frac{1}{\mu_X}$, which

would naturally be estimated by $g(\bar{X}_n) = \frac{1}{\bar{X}_n}$.

$$E(\bar{X}_n) = \mu_X$$

$$V(\bar{X}_n) = \frac{\sigma^2}{n}$$

$$g(x) = \frac{1}{x} = x^{-1}$$

$$g'(x) = -\frac{1}{x^2}$$

$$g'(\mu_X) = -\frac{1}{\mu_X^2}$$

$\bar{X}_n \sim \text{Normal}$
by CLT

so Δ -method says $g(\bar{X}_n) = \frac{1}{\bar{X}_n} \sim \text{Normal}$

with mean $g(\mu_X) = \frac{1}{\mu_X}$ and variance

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

$$\sigma_X^2 / (n \mu_X^4)$$

Specific
Calculation

Under some plausible assumptions, (318)
we've seen that $(X_i | \lambda) \stackrel{\text{IID}}{\sim} \text{Exponential}(\lambda)$

may be a reasonable model for waiting times.

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

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

so $\frac{1}{\bar{X}_n}$ should (large n)

be approximately Normal with mean $\frac{1}{\lambda} = \lambda^{-1}$

$$\text{and SD } \frac{\sigma_{\bar{X}}}{\mu_{\bar{X}} \sqrt{n}} = \frac{\frac{1}{\lambda}}{\left(\frac{1}{\lambda}\right)^2 \sqrt{n}} = \frac{\lambda}{\sqrt{n}}$$

(discrete or
continuous)

Fourier version
of Δ -method

$\mathcal{F}_1, \mathcal{F}_2, \dots$ sequence of v.v.;
 F^* continuous cdf;

θ a real number; $a_1, a_2, \dots \uparrow \infty$
positive sequence

$g(\cdot)$ a ^{real-valued} function of a real variable (319)
 such that $g'(\cdot)$ is continuous and
 $g'(\theta) \neq 0$; then if $a_n(\bar{Y}_n - \theta) \xrightarrow{D} F^*$,

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

Typical application:
 X_1, X_2, \dots IID

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

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

In this context the theorem says that

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

(28 Aug 17)
~~(27 Aug 17)~~

is also $\sim N(0, 1)$

A little bit more about the continuity correction

T97-Suchs case study, revisited

$$X = \# \text{ T-S babies}$$

in family of $n=5$ children, both parents carriers so that

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

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

Then $(T_i) \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p)$ and $X = \sum_{i=1}^n T_i$
($i=1, \dots, n$) $i=1, \dots, n=5$

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

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

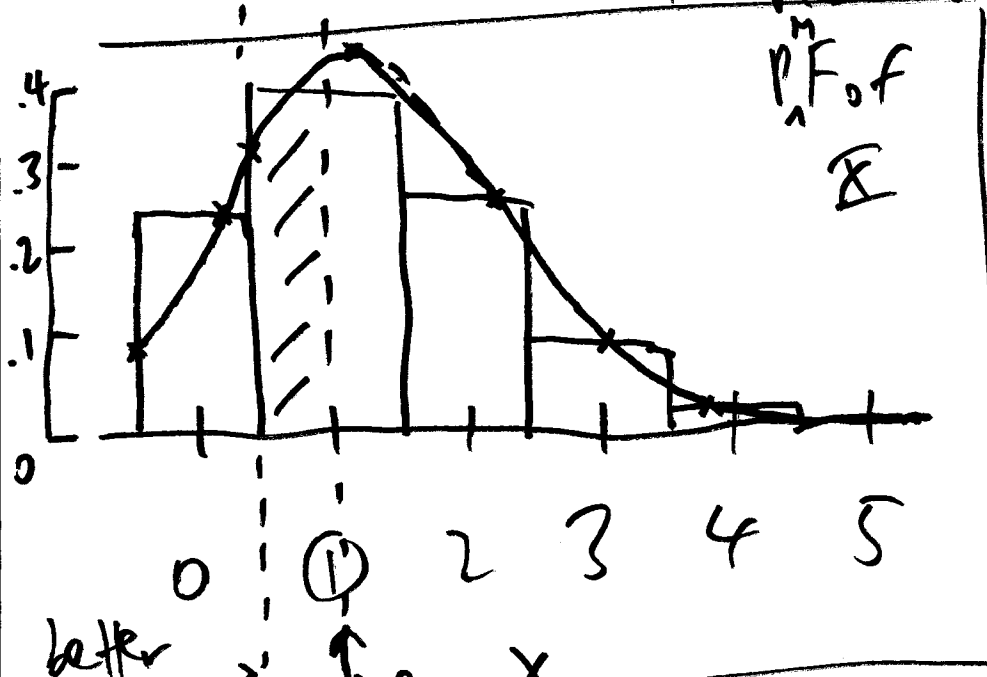
$$\sigma_{\bar{X}} = \sqrt{V(\bar{X})} = \sqrt{np(1-p)} = 0.98 \quad (32)$$

on day 1 of this class we worked out

that $P(\text{1 or more T-S babies}) = P(\bar{X} \geq 1)$

$$1 - P(\text{no T-S babies}) = 1 - (1-p)^n = 0.76$$

$$= 1 - P(\bar{X} = 0)$$

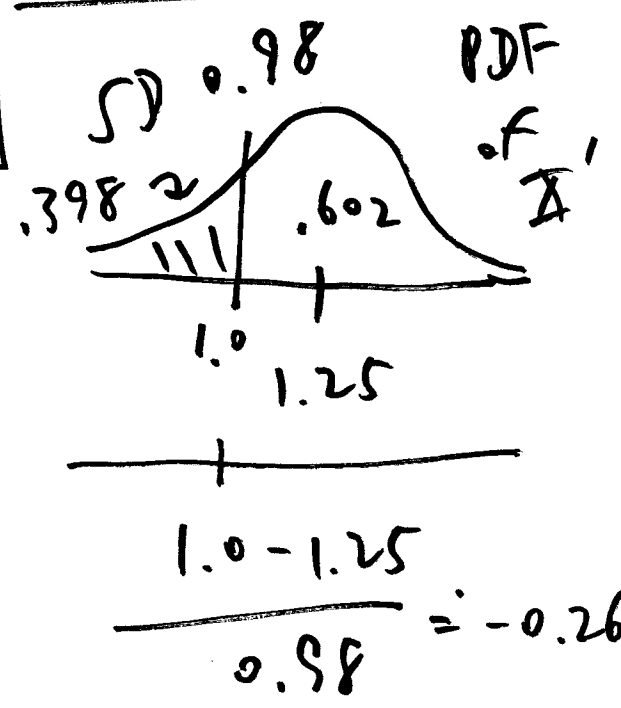


Naive Normal approximation, from CLT:

$$P(\bar{X} \geq 1) = 1 - P(\bar{X}' < 1)$$

$$= 1 - 0.398$$

$= 0.602$ (quite a bad approximation)



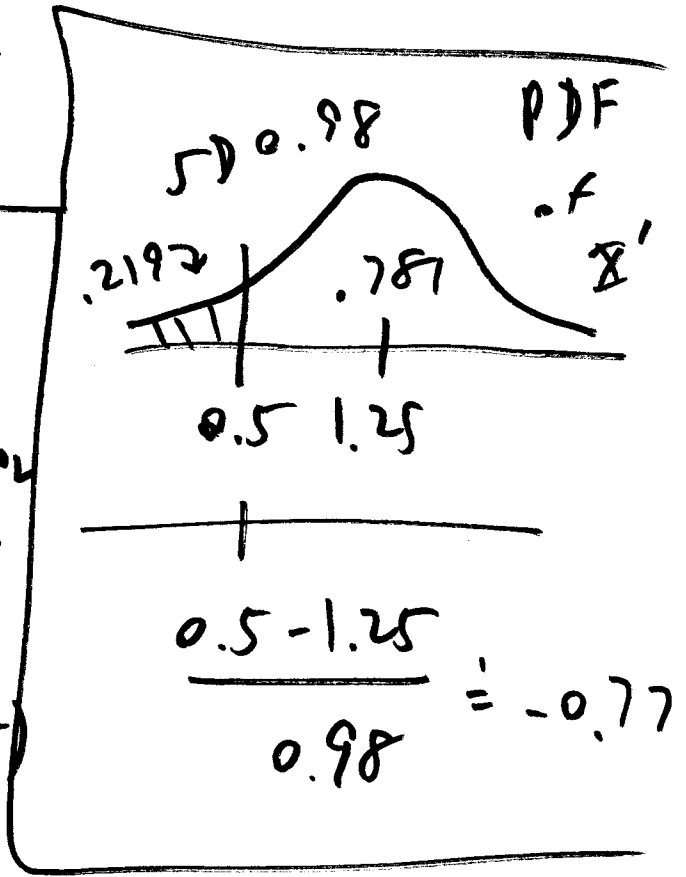
Improved approximation obtained by paying attention to the edges of the histogram ($\frac{M}{n}$) bars:

Normal approximation with continuity correction

$$\begin{aligned}
 P(X \geq 1) &= 1 - P(X' < 0.5) \\
 &= 1 - .219 \\
 &= 0.781
 \end{aligned}$$

(correct answer 0.76; much better approx.)

(4 Jan 19)



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

Markov Chains

Recall the definition of a stochastic process:

Def. A sequence of rvs X_1, X_2, \dots (323)
is called a stochastic process with
discrete time parameter $t = 1, 2, \dots$.

X_1 is the initial state of the process;

$X_n, n \geq 1$ is the state of the process
at time $t = n$.

The simplest possible
discrete-time stochastic process is
an IID sequence of rvs (X_1, X_2, \dots) .

Suppose that there's a parameter θ
such that $(X_i | \theta) \stackrel{\text{IID}}{\sim}$ from some dist.

depending on θ .

Q: Does this process
have a memory?

Example,
revisited

Machine with a dial from θ (324)
0 to 1, produces IID Bernoulli(θ)

Recall that
trials X_i : The process (X_1, X_2, \dots)

does not have a memory ^{for you} if θ is unknown

to you: the information that 17 out of the first 20 trials were successes helps you to predict X_{21} , because it's reasonable to conclude from X_1, \dots, X_{20} that θ is around $\frac{17}{20} = 0.85$, so X_{21} ~~is~~ ^{will}

probably ^{be} a success.

But the process

$\{(X_i | \theta), i=1, 2, \dots\}$ has no memory

once θ is known: information about

The first n trials is irrelevant to (325)
your prediction of X_{n+1} if you know

θ . An IID process $(X_i | \theta) \stackrel{\text{IID}}{\sim}$
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}, \dots
(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 θ in the notation.

discrete-time

Def.

A stochastic process is a

(first-order) Markov chain if for $n = 1, 2, \dots$; b any real number; and for all possible sequences of states x_1, x_2, \dots

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

In other words, the only thing you need to know to simulate where the Markov chain is going next is where it is now. (28 Aug 19)

(Can define higher-order Markov chains 327 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. (countably or uncountably)

(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 P^N of a Markov chain with finite S :

Consequences

Def. A Markov chain with a finite state space is called a finite Markov chain.

① (X_1, X_2, \dots) finite

328

Markov chain \rightarrow

$$P(X_1 = x_1, \dots, X_n = x_n) =$$

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

$$P(X_3 = x_3 | X_2 = x_2) \cdot \dots$$

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

Suppose you have a

Def. finite Markov chain with k possible states numbered $1, \dots, k$

(k integer ≥ 2) $\rightarrow \{P(X_{n+1} = j | X_n = i),$

$i, j = 1, \dots, k, n = 1, 2, \dots\}$ ~~is~~ called the transition

distribution of the Markov chain.

If $P(X_{n+1}=j | X_n=i)$ is the same

for all n , the transition distribution

is said to be stationary ^{(DS) (bed name)}. If (time-homogeneous)

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

$P_{ij} \triangleq P(X_{n+1}=j | X_n=i)$ completely

characterize the Markov chain's

behavior.

in a matrix called the transition matrix.

Can arrange the P_{ij} to state P_{ij}

$$\begin{matrix}
 & \begin{matrix} 1 & 2 & \dots & k \end{matrix} \\
 \begin{matrix} 1 \\ 2 \\ \vdots \\ k \end{matrix} & \begin{bmatrix} p_{11} & p_{12} & \dots & p_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ p_{k1} & p_{k2} & \dots & p_{kk} \end{bmatrix}
 \end{matrix}$$

from state k

All of the elements of \underline{P} are (330)
(and ≤ 1)
non-negative (they're probabilities), and
all of the row sums are 1 (because
the chain has to go somewhere), i.e.

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

matrix versus quaternion

A square matrix \underline{P} with non-negative
entries and ^{all} row sums equal to 1
is called a stochastic matrix.

~~(Definitive)~~
Example [^] Gene inheritance is Markovian:
genetic makeup at birth
you, is the genetic story of your parents

(your grand parents, ..., are irrelevant) (33)

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

{ allele 1 from parent 1, allele 2 from parent 1, allele 1 from parent 2, allele 2 from parent 2 }, for

example {Aa, Aa}.

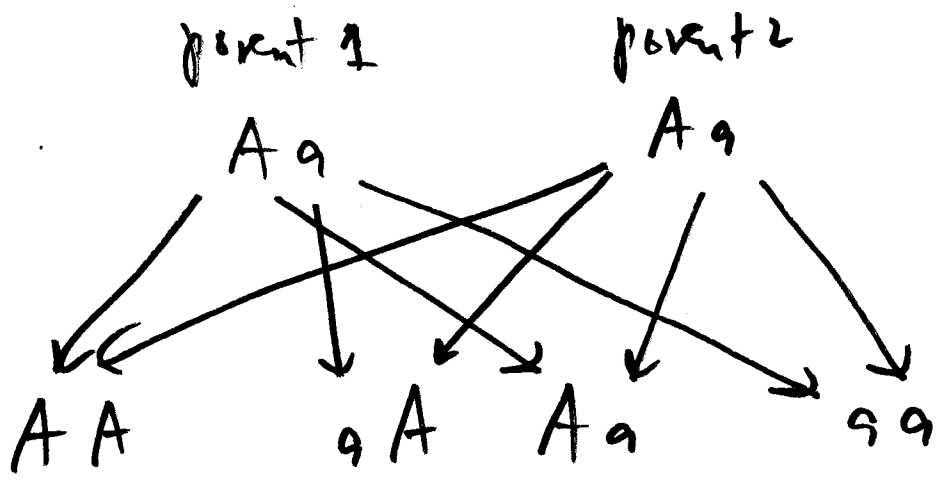
Ignoring order

(because it's irrelevant in inheritance),

there are 6 possible states: {AA, AA}

{AA, Aa}, {AA, aa}, {Aa, Aa}, {Aa, aa}

and {aa, aa}.

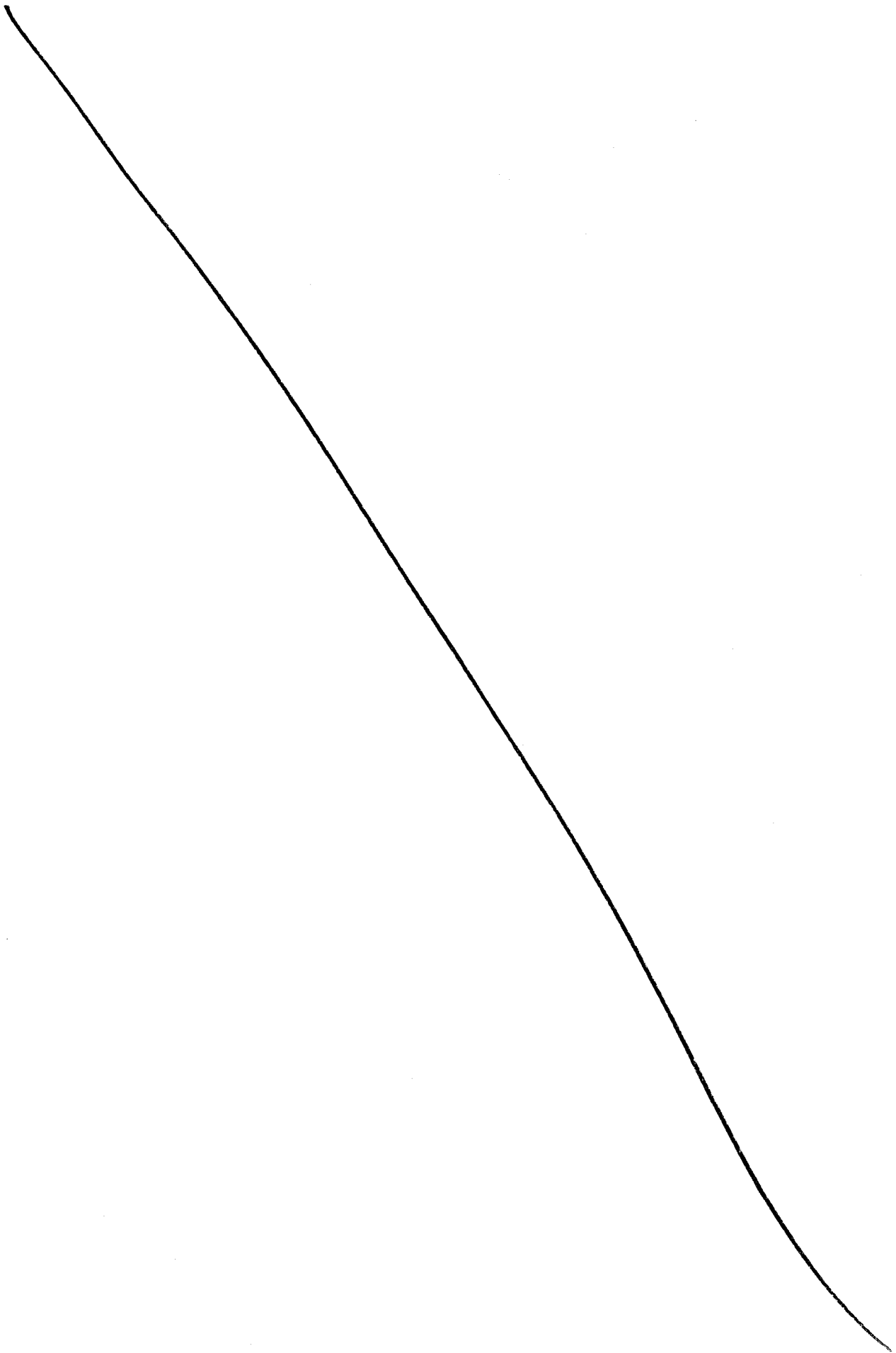


one possible inheritance sequence

offspring gets A or a from parent 1 and A or a (independently) from parent 2, 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	0	0	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$
{aa, aa}	0	0	0	0	0	1



Example (random walk) You're watching 334

a particle move around on the

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

over time: here are the rules:

whenever 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?

	to → ...	-2	-1	0	1	2	...	
from ↓	⋮	⋮	⋮	⋮	⋮	⋮	⋮	
-2	...	p_2	p_3	0	0	0	...	
-1	...	p_1	p_2	p_3	0	0	...	
0	...	0	p_1	p_2	p_3	0	...	$= P$
1	...	0	0	p_1	p_2	p_3	...	
2	...	0	0	0	p_1	p_2	...	
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	

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, P is said to be tridiagonal.

Moreover, all of the main diagonal entries are the same (p_2); all of the entries 1 diagonal ~~above~~ ^{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)

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

(337)

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 $n \rightarrow \infty$ 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.

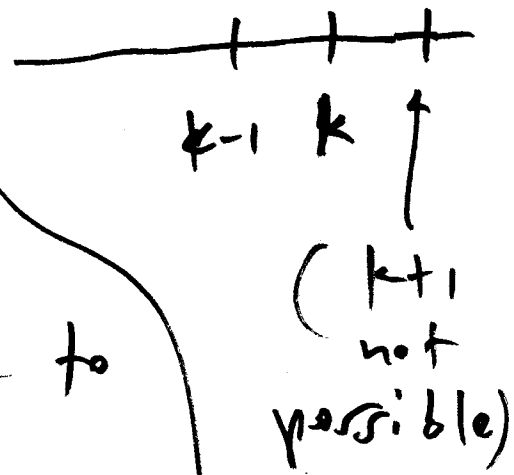
The infinite random walk evidently has "too much freedom" to move around to get interesting results; let's bound it.

Bounded
random
walk

Restrict the Markov chain (338)
to $S = \{-k, -k+1, \dots, -1, 0, 1, \dots, k-1, k\}$
for some integer $k \geq 1$.

Q: what
to do at the boundaries?

one idea would be to wrap
around: if you try to move to
($k+1$), interpret that as a move
to $-k$; if you try to move to $-(k+1)$,
move to $+k$.



(move.2)

Transition matrix with $k=2$

R demo

Another idea:
keep trying
until you make
a legal move (move.1)

from

	-2	-1	0	1	2	$M(k)$
-2	p_2	p_3	0	0	p_1	
-1	p_1	p_2	p_3	0	0	
0	0	p_1	p_2	p_3	0	
1	0	0	p_1	p_2	p_3	
2	p_3	0	0	p_1	p_2	

(move.2)

Back to
a general
finite
Markov
chain

Let $p_{ij}^{(m)} = P(\text{chain moves from } (i) \text{ to } (j) \text{ in } m \text{ steps})$ (339)

Theorem

$$= P(X_{n+m} = j \mid X_n = i)$$

Finite Markov chain with stationary transition distributions & transition

matrix $\underline{P} \rightarrow p_{ij}^{(m)}$ is just the (i, j)

entry of the matrix \underline{P}^m , which

is called the m -step transition matrix

of the Markov chain.

Genetic example,
continued

$\{AA, AA\}$ has the property that once the chain is in that state, it can't

so anywhere else; so does $\{aa, aa\}$ (34)

This occurs for a state i when $p_{ii} = 1$.

Def. Any state with $p_{ii} = 1$ is

called an absorbing state.

Notice

that in this genetic Markov chain,
states ~~1, 2~~ ^{1, 2, 4} all have positive probability
of moving to state 1 in 2 steps,
and the same is true of moving to
state 6 in 2 steps.

It follows that,

if the chain is run long enough (simulating
many generations), it will either end up

in state $\{AA, AA\}$ or in state $\{aa, aa\}$ (341)

Markov chains that settle down to a single ^{stable} long-run distribution are especially important in contemporary Bayesian computation; the long-run stable distribution is called the equilibrium distribution of the chain.

Important
note on
terminology

JS call this distribution the stationary dist. of the chain, but this choice is unfortunate because they've already used stationary to mean something else:

DS: (If $P(X_{n+1}=j | X_n=i)$ is the same for all n , DS say) that the transition distribution is stationary; other people call this time-homogeneous.

I'll use equilibrium distribution for

the long-run behavior of Markov chains that settle down into a stable long-run story.

Where should the Markov chain start?

You can either initialize a Markov chain to a deterministic value, or you can start it off by making a ^{random} draw from what's called the initial distribution of the Markov chain:

Def Any vector \vec{v} of non-negative numbers that add up to 1 is called a probability vector; any such vector whose components specify that a Markov chain will be in each possible state at time 1 is referred to as the initial distribution of the chain.

So: After 1 timestep, ^(iteration) the probability dist. over the Markov chain's possible states is \vec{v} ; after 2 iterations the chain's dist. is $\vec{v} P$; after $(n+1)$ iterations its dist. is $\vec{v} P^n$; it would be nice if $\vec{v} P^n$ converged to a unique dist. as $n \rightarrow \infty$: this would

be its equilibrium distribution, (344)

Notice something interesting: if we choose \underline{v} so that $\underline{v} \underline{P} = \underline{v}$, then

$$\underline{v} \underline{P}^2 = \underline{v} \underline{P} = \underline{v}, \quad \underline{v} \underline{P}^3 = (\underline{v} \underline{P}^2) \underline{P} = \underline{v} \underline{P} = \underline{v}$$

$$= \underline{v} \underline{P} = \underline{v}; \quad \text{and so } \lim_{n \rightarrow \infty} \underline{v} \underline{P}^n = \underline{v}$$

Def. Markov chain with transition matrix \underline{P} → any probability vector \underline{v} such that $\underline{v} \underline{P} = \underline{v}$ is an equilibrium dist. for the Markov chain

under additional

conditions on \underline{P} , such an equilibrium dist. will be unique (we won't fully pursue that here).

How find \underline{v} so that $\underline{v} \underline{P} = \underline{v}$? (345)

In linear algebra this is an example of an eigenvalue/eigenvector problem:

Def. Given a square matrix $\underline{P}_{k \times k}$,

any vector $\underline{v}_R \in \mathbb{R}^k$ satisfying $\underline{P}_{k \times k} \underline{v}_R = \lambda_R \underline{v}_R$

is called a right eigenvector of \underline{P} with

(right) eigenvalue λ_R , and any vector $\underline{v}_L \in \mathbb{R}^k$

satisfying $\underline{v}_L \underline{P}_{k \times k} = \lambda_L \underline{v}_L$ is called

a left eigenvector of \underline{P} with (left)

eigenvalue λ_L .

So, given a transition matrix P for a Markov chain, an equilibrium dist. for the chain can be found by computing the left eigenvector \underline{v}_k where eigenvalue is 1, if such a vector exists.

Most computer routines $\underline{v}_k P = \underline{v}_k$

for eigenanalysis only give you right eigenvectors, but notice that if

$$\underline{v}_k P = \underline{v}_k \text{ then } \left(\underline{v}_k P \right)^T = \underline{v}_k^T$$

← transpose

$$P^T \underline{v}_k^T = \underline{v}_k^T \text{ so we can just}$$

eigendecompose P^T instead of P .

Genetic
example,
continued

R 's routine eigen gives (347)
the following results: \underline{p}^T has

two eigenvectors whose
eigenvalues are 1: $\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$ and $\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$,
corresponding to the
two absorbing states.

This suggests that there's an entire family
of equilibrium distributions, of the

form $(p, 0, 0, 0, 0, 1-p)^T$ for
 $0 \leq p \leq 1$; and \boxed{Wd} can be used to
~~check~~ verify this conjecture.

So the earlier guess is also correct:

after many generations either one of
 $\{AA, AA\}$ or $\{aa, aa\}$ will be absorbing.

There is a special case in which a unique stationary distribution exists.

Theorem If you can find a positive integer $n \geq 1$ such that every element of P^n is strictly positive, then $\lim_{n \rightarrow \infty} P^n$ is a matrix with all rows equal to the unique stationary dist \underline{v} , and no matter what the chain's initial distribution is, its distribution after n steps converges to \underline{v} as $n \rightarrow \infty$.
