

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 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⁾
existing
→ for every $t \geq 0$,

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

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 for most distributions it's (also) crude (as with the Markov bound). Back to \bar{X}_n

$X_i \stackrel{iid}{\sim}$ 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 $V(\bar{X}_n) = \frac{\sigma^2}{n}$.

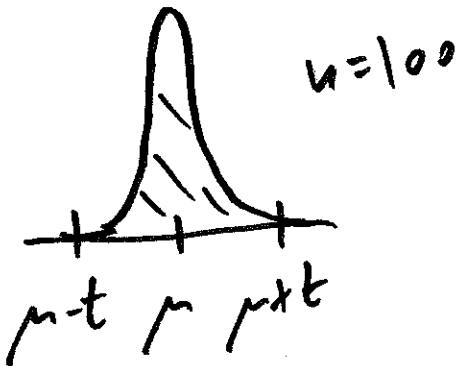
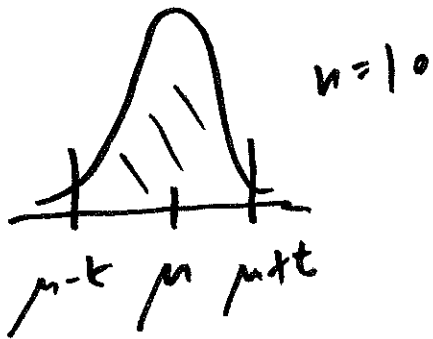
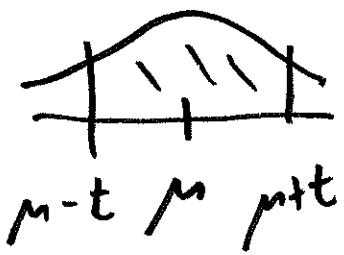
Chebyshev then

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

this can be

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

PDF of \bar{X}_n $n=1$



⋮

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 $Z_n \xrightarrow{P} b$.

An immediate

consequence of Chebyshev & this definition is

(weak)
Law of
Large
Numbers

$X_i \stackrel{IID}{\sim}$ 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 (1655-1705) proved it for $(X_i | \theta) \stackrel{IID}{\sim}$ Bernoulli(θ) (it took him 20 years to find ~~the~~ 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 $g(Z_n) \xrightarrow{P} g(b)$.

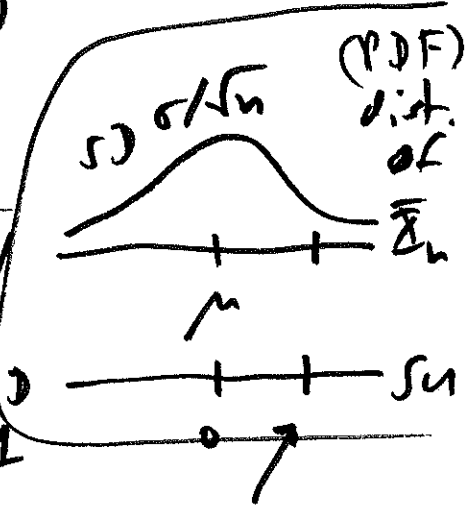
Central Limit Theorem (CLT)

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

$X_i \sim ?$

A: Yes: it's the most famous result in all of probability:
 $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$

Central Limit Theorem (CLT)

$X_i \sim \text{any}$ 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 ^(3.0)
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 $0 < \sigma^2 < \infty$, $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

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

the
CLT

also has a long history: it was

first demonstrated for $X_i \sim \text{Bernoulli}(p)$ ^{IID}
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. CLT name due to
Hungarian-American mathematician (1887-1985) George Pólya in
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.

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 ^{desired} accuracy target? (310)

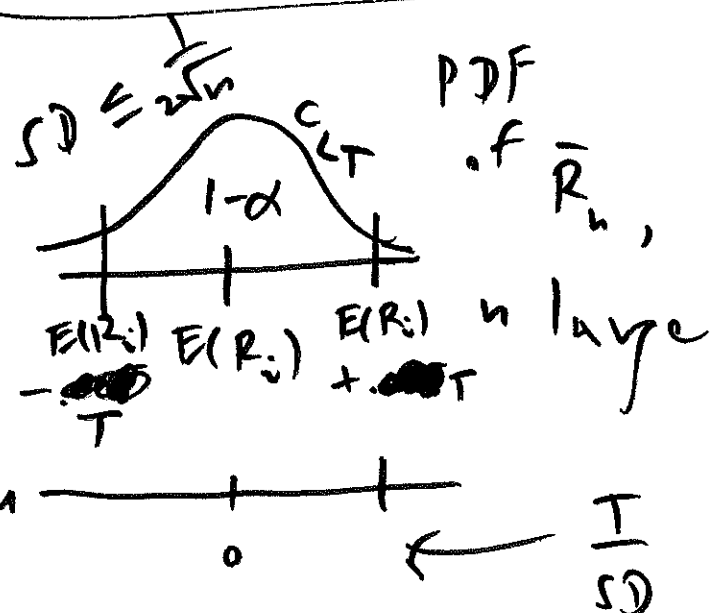
By definition

$$|R_i| = \left| \frac{Y_i}{\Sigma_i + Y_i} \right| \leq 1; \text{ can show that}$$

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

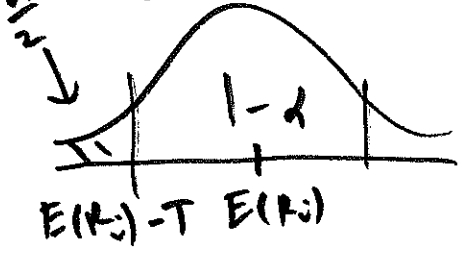
Says that dist. of \bar{R}_n will be close to Normal for large n , with mean $E(R_i)$

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



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

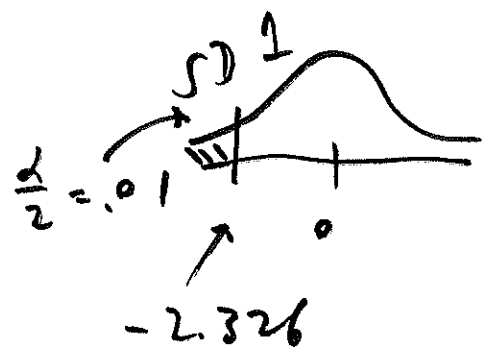
$$\Phi^{-1}\left(\frac{\alpha}{2}\right) = \frac{[E(R_i) - T] - E(R_i)}{SD} = \frac{-T}{SD} \leq 2T\sqrt{n}$$

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 $\alpha = .02$ to get

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



simulation replications

needed

Case Study: Escalators

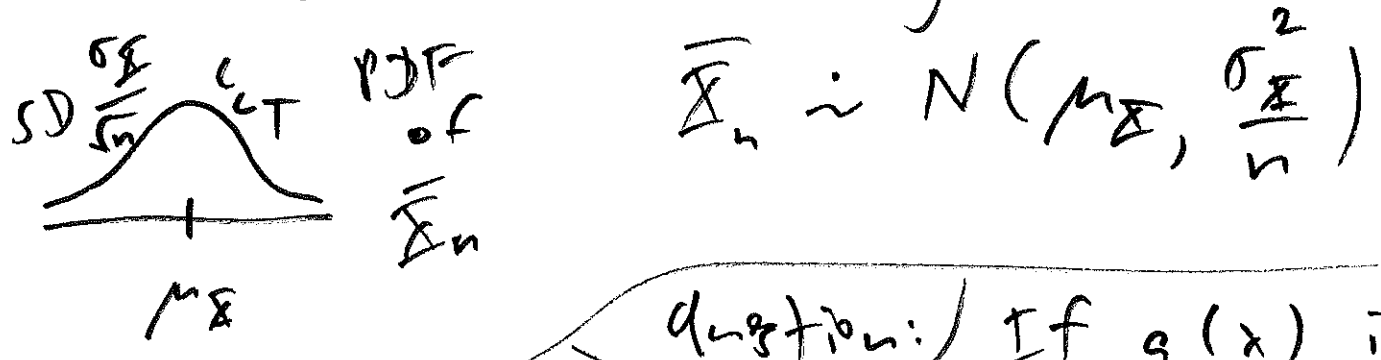
in the London Underground (👤)

The Delta Method

The CLT says that if $X_i \stackrel{iid}{\sim}$ (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



Question: If $g(x)$ is

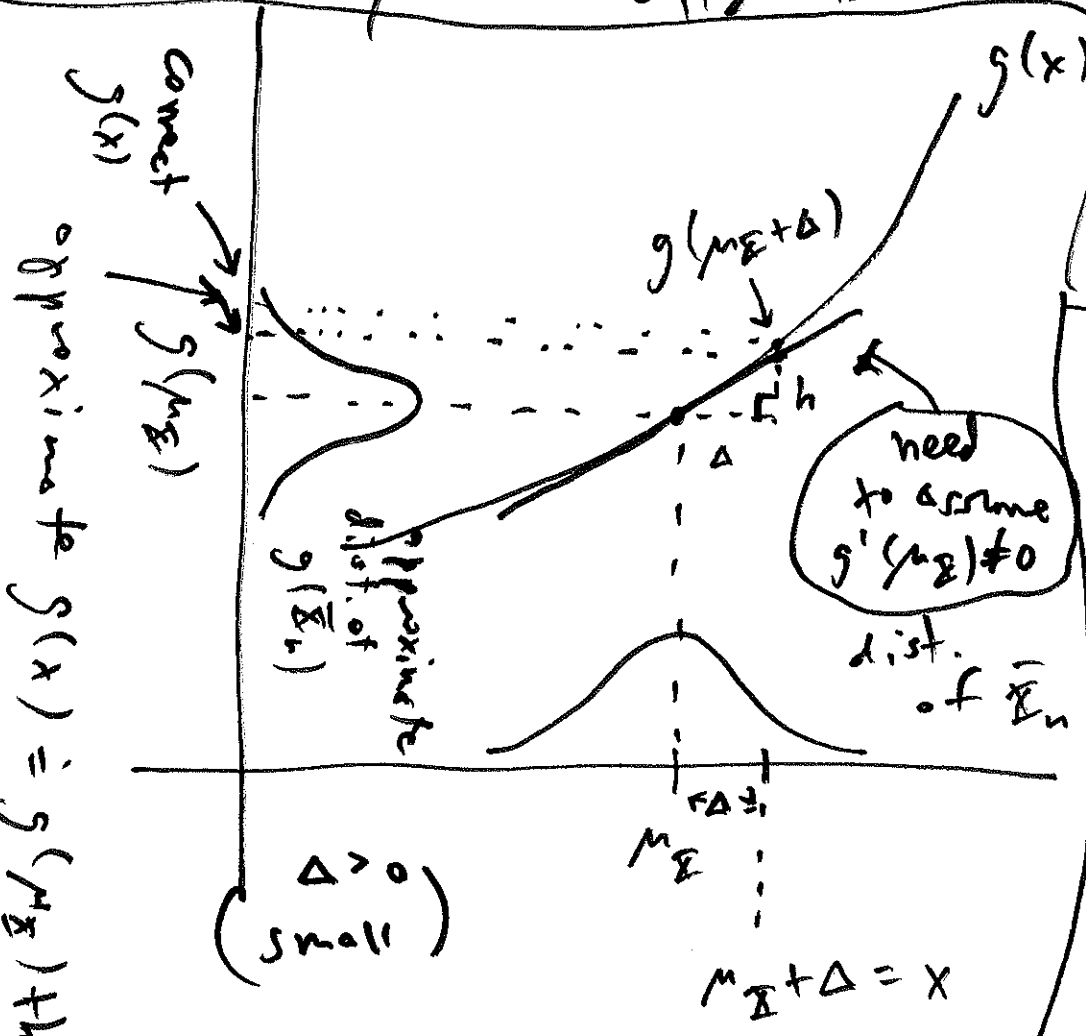
a sufficiently "nice" function, is there a comparable result for $g(\bar{X}_n)$?

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

\bar{X}_n should be close to $\mu_{\mathbb{R}}$ 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

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

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



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

so

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

$$= g(\mu_{\mathbb{R}}) + g'(\mu_{\mathbb{R}}) \cdot \Delta$$

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

so $\Delta = x - \mu_{\mathbb{R}}$

$$g(\bar{X}_n) = g(\mu_X) + g'(\mu_X)(\bar{X}_n - \mu_X) \quad \text{so} \quad (314)$$

\uparrow constant $\quad \uparrow$ r.v. $\quad \uparrow$ r.v.

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

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

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

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

\uparrow constant $\quad \uparrow$ r.v.

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

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

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

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 :

\forall 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} $E(X_i) = \mu_X$

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

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

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

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

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

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.

$E(X_i) = \frac{1}{\lambda}$, $V(X_i) = \frac{1}{\lambda^2}$ $(X_i | \lambda)$ has PDF
 $= \mu_X$, $= \sigma_X^2$
 $f_{X_i}(x_i | \lambda) = \lambda e^{-\lambda x_i} I(x_i > 0)$

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

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

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

(discrete or continuous)

Fancy version
of Δ -method

$\mathcal{I}_1, \mathcal{I}_2, \dots$ sequence of i.i.d.
 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 $a_n(\bar{Y}_n - \theta) \xrightarrow{D} F^*$,

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

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)

~~(29 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 $(T_i) \stackrel{i.i.d.}{\sim} \text{Bernoulli}(p)$ and $X = \sum_{i=1}^n T_i$
($i=1, \dots, n$)

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

$$\mu_X = E(X) = np = 1.25 \text{ and } 5)$$

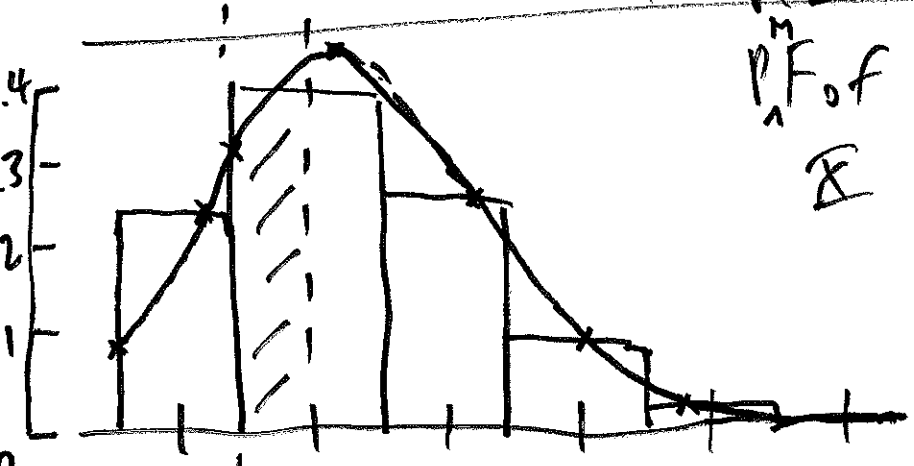
$$\sigma_{\bar{X}} = \sqrt{V(\bar{X})} = \sqrt{np(1-p)} \approx 0.98 \quad (32)$$

on day 1 of this class we worked out

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

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

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



PDF of \bar{X}

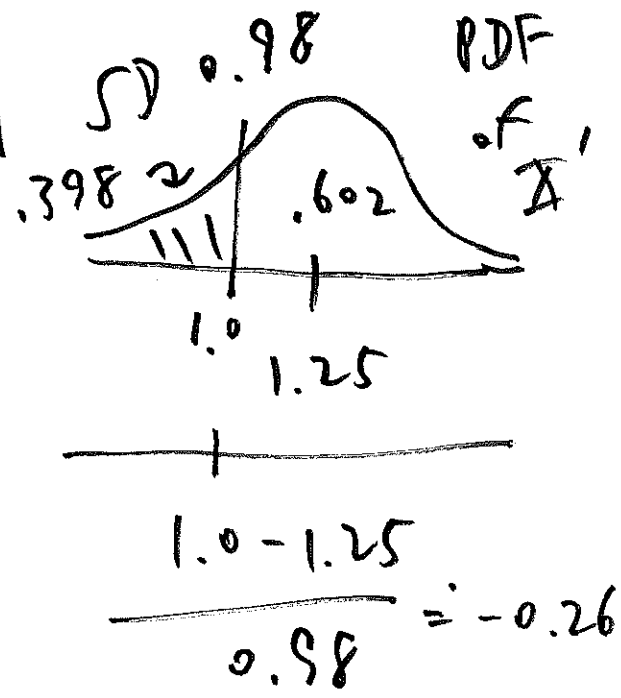
Naive Normal approximation, from CLT:

0 1 2 3 4 5
 better approx → naive approx

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

$$= 1 - 0.398$$

≈ 0.602 (quite a bad approximation)



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

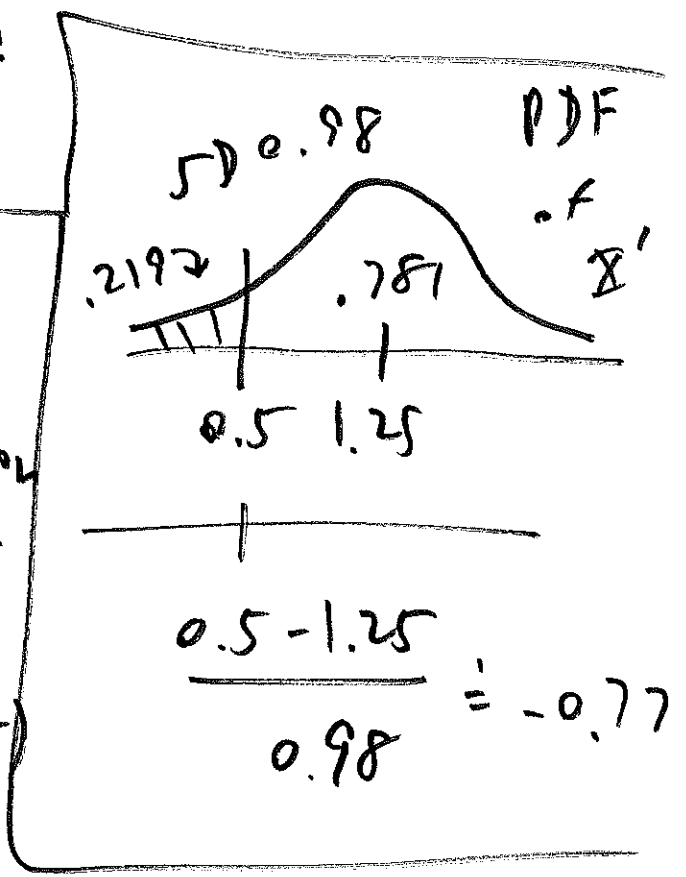
Normal approximation with continuity correction

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

$$= 1 - .2192$$

$$= 0.781$$

(correct answer 0.76; much better approx.)
(4 Jan 19)



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

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

does ^{for you} have a memory 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

Q. 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).