

and $P(X=x | Y=y) = \binom{y}{x} p^x (1-p)^{y-x}$ (25)

Notice that if $X=x$, $Y \geq x$ because the ^{actual} number of oocysts (Y) has to be at least as large as the number of oocysts detected (X).

After a careful

$$f_X(x) = \sum_{y=x}^{\infty} \binom{y}{x} p^x (1-p)^{y-x} \frac{(\lambda t)^y e^{-\lambda t}}{y!}$$

calculations you get;

$$= \frac{e^{-p\lambda t} (p\lambda t)^x}{x!}$$

i.e., $X \sim \text{Poisson}(p\lambda t)$:

losing a proportion

$(1-p)$ of the oocysts to faulty counting

just lowers the rate of the Poisson

process from λ /liter to $\lambda \cdot p$ /liter

(makes excellent sense).

Negative Binomial Distribution

You're watching a potential ⁽²⁵³⁾ endless sequence of Bernoulli trials with constant success

probability p .

Let X = # failures before $(r)^{\text{th}}$

You can show that X ^{what's called}

follows the Negative Binomial dist:

its PF is $f(x | r, p) = \binom{r+x-1}{x} p^r (1-p)^x$.

with parameters (r, p)

The name comes ^($0 < p < 1$) $X \in \{0, 1, 2, \dots\}$ $\binom{r+x-1}{x}$.

from the fact that, when you watch a sequence of Bernoulli trials with constant ^{unknown} success probability p unfold, there are two different ways to

estimate p : decide ahead of time to (254)
(known constant)
sample n success/failure trials, and
record the (random) # S of successes
you see (from which a reasonable
estimate would be $\hat{p}_B = \frac{S}{n}$ ← Binomial)

(or) decide ahead of time that you're
going to sample until you've seen s
(known constant) successes & record the
(random) # of trials N needed
to accumulate that many successes
(from which a reasonable estimate
would be $\hat{p}_{NB} = \frac{s}{N}$ ← Negative Binomial).

Special
Case of
Negative
Binomial

Set $r=1$ and record the (255)
number X of failures until
the first success: X is
said to follow the

Geometric (p) distribution, with

$$P\{X=x\} = p(1-p)^x \mathbb{1}_{\{0,1,\dots\}}(x)$$

(parameter p)

~~Con~~ sequence X_1, \dots, X_n IID Geometric(p)

$$\sum_{i=1}^n X_i \sim \text{Negative Binomial}(n, p)$$

This is a direct analogue to the

Bernoulli/Binomial story: X_1, \dots, X_n IID

$$\text{Bernoulli}(p) \rightarrow \sum_{i=1}^n X_i \sim \text{Binomial}(n, p)$$

$X \sim \text{Negative Binomial}(r, p)$

256

$$\psi_X(t) = \left[\frac{p}{1 - (1-p)e^t} \right]^r \quad \text{for } t < \log\left(\frac{1}{1-p}\right)$$

from which $E(X) = \frac{r(1-p)}{p}$, $V(X) = \frac{r(1-p)}{p^2}$

Consequence

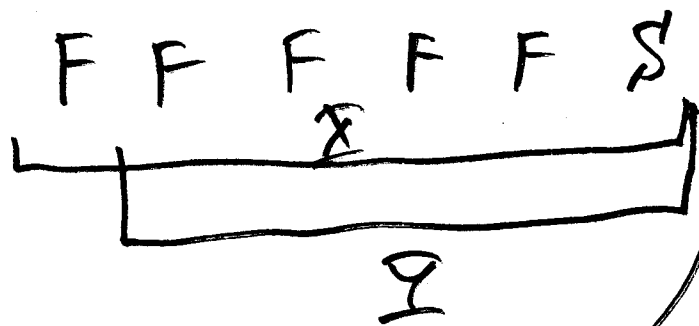
$X \sim \text{Geometric}(p) \rightarrow$

$\begin{cases} k \\ t \end{cases}$ both non-negative integers

$$P(X = k+t \mid X \geq k) = P(X = t)$$

this is called the memoryless property of the Geometric distribution, and it turns out that this is the only

discrete distribution with this property (257)



$X = \#$ failures until first success = 5 (here)

$Y = \#$ failures, starting at trial $(k+1)$ until next success (here, = 2)
 (- 4 here) Then Y has

the same dist. as X and is independent of what happened on the first k trials, i.e., "the process has no memory".

Core 2: Important Continuous Distributions

Normal (Gaussian) Distribution

$X \sim \text{Normal}(\mu, \sigma^2)$ mean μ variance $\sigma^2 < \infty$

PDF \rightarrow

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$$

The Normal dist. is the single most important dist. in all of probability & statistics, mainly for 2 reasons:

- ① many observable random processes have dist. shapes that are close to the "bell curve" (Normal PDF), and
- ② the Central Limit

Theorem (CLT), which we'll examine soon.

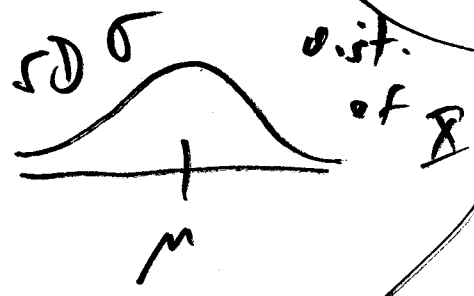
Properties of the Normal Dist.

$N(\mu, \sigma^2)$

$X \sim \text{Normal}(\mu, \sigma^2) \quad | \quad E(X) = \mu$

$V(X) = \sigma^2, \quad SD(X) = \sigma$

$\Psi_X(t) = \exp\left(\mu t + \frac{\sigma^2 t^2}{2}\right)$



(center of symmetry)
mean
median
mode
= μ

Consequences ① $X \sim \text{Normal}(\mu, \sigma^2)$, (259)

$Y = aX + b$, ($a \neq 0$) fixed constants \rightarrow

$Y \sim \text{Normal}(a\mu + b, a^2\sigma^2)$.

In other words, Normality is preserved under linear transformations Def.

The Normal dist. with mean $\mu = 0$ and SD $\sigma = 1$ is called the standard normal dist.

The PDF of $X \sim \text{Normal}(0, 1)$ is

$\phi_X(x) \triangleq \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$ and its
 ϕ (lower-case)

CDF is $F(x) \triangleq \int_{-\infty}^x \phi_X(t) dt$
 F (upper-case)

It turns out that e^{-cx^2} has no (260)
anti-derivative in closed form, so
 $\Phi(x)$ cannot be summarized in a
formula; instead it's approximated by
numerical integration (see p. 861 in DS).

Consequences,
continued

② Because the Normal PDF
(for all $x \in \mathbb{R}$)
is symmetric, $\Phi(-x) = 1 - \Phi(x)$

$$\text{and } \Phi^{-1}(p) = -\Phi^{-1}(1-p) \text{ (for all } 0 < p < 1)$$

③ $X \sim \text{Normal}(\mu, \sigma^2) \rightarrow Z = \frac{X - \mu}{\sigma} \sim N(0, 1)$

$$\text{so that } F_X(x) = \Phi\left(\frac{x - \mu}{\sigma}\right)$$

$$\text{and } F_X^{-1}(p) = \mu + \sigma \Phi^{-1}(p)$$

Empirical Rule

Part 1 Start at the mean μ (261) of a distribution and go $\pm 1\sigma$

either way: you will find (about $\frac{2}{3}$) (68%) of the probability in the

interval $(\mu \pm 1\sigma)$ Part 2 Ditto 2SDs

either way: $(\mu \pm 2\sigma)$ captures (about ^{most} 95%) of the probability

Part 3

Ditto 3SDs either way: $(\mu \pm 3\sigma)$ captures almost all (99.7%) of the

probability

This Rule is exact for all Normal dists & is a surprisingly

good approximation for many other (262)

distributions.

This permits an easy trick

that's helpful in computing Normal probabilities.

Example:

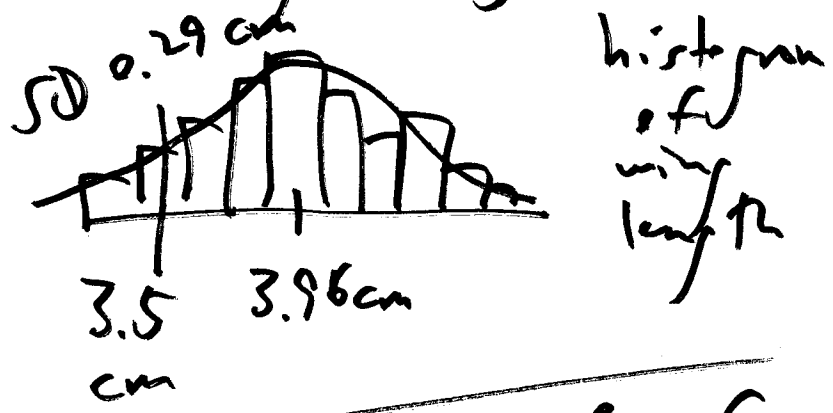
You have a random sample

of $n = 103$ immature monarch butterflies, and you measure their wing lengths:

$y = \text{wing length (cm)}$

$y_1 = 4.1$
$y_2 = 3.3$
\vdots
$y_n = 4.7$

$n = 103$



mean $\bar{y} = 3.96 \text{ cm}$
 SD $s = 0.29 \text{ cm}$

Q: About what % of the sampled butterflies had wing length $\leq 3.5 \text{ cm}$?

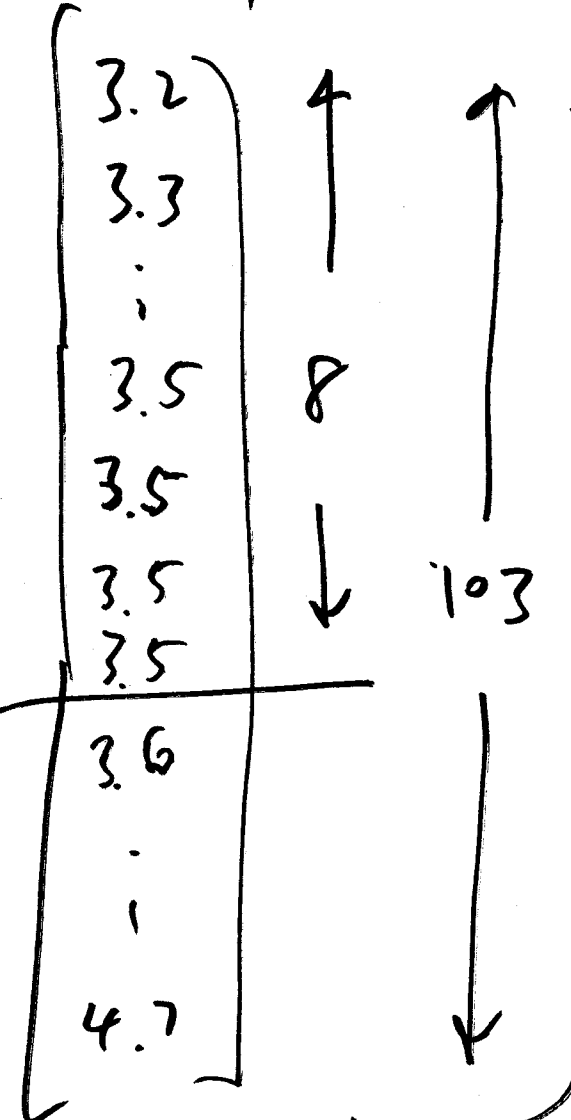
$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

sample mean

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$$

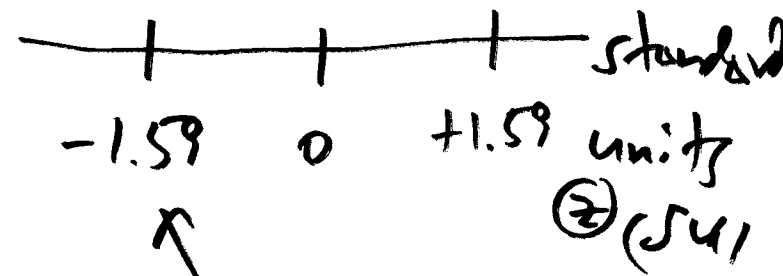
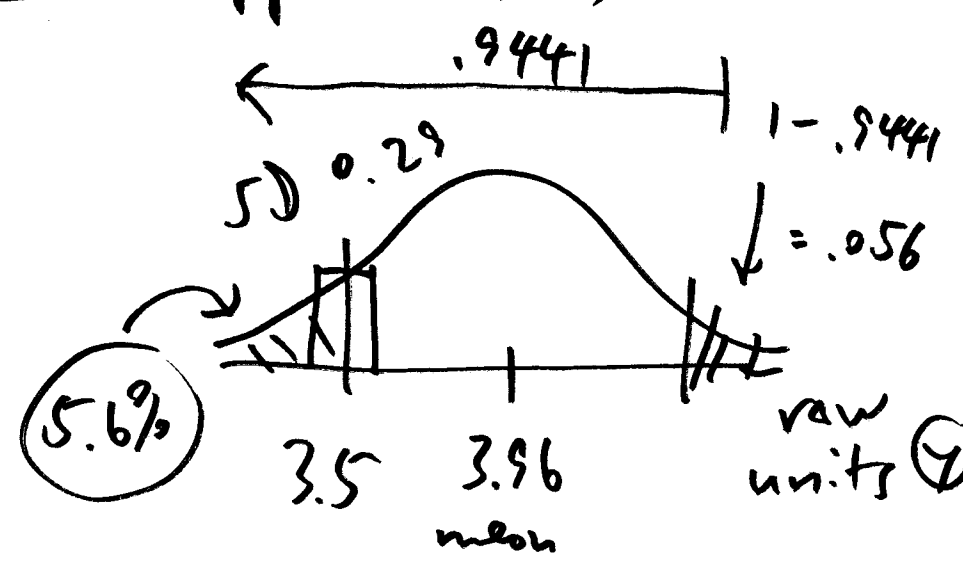
sample SD

sorted y



A_1 (exact) $\frac{8}{103} = 7.8\%$ (263)

A_2 (approximate)

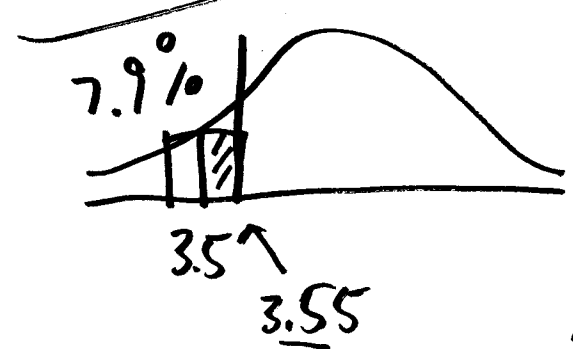


continuity to f_u
for data:

$$z = \frac{y - \bar{y}}{s} = 5u$$

for random variables

$$z = \frac{Y - \mu}{\sigma} = 5u$$



keeping track of histogram bar edges: continuity correction

More consequences

(4) X_1, \dots, X_k independent,

$X_i \sim \text{Normal}(\mu_i, \sigma_i^2)$

$\rightarrow \sum_{i=1}^k X_i \sim \text{Normal}(\sum_{i=1}^k \mu_i, \sum_{i=1}^k \sigma_i^2)$

nice additive property

this is why Normal dists are indexed by variance rather than SD.

Notation

$\text{Normal}(\mu, \sigma^2) \triangleq N(\mu, \sigma^2)$

Example Population of ^{adult u.s.} women: height follows $N(\mu = 65.0 \text{ in}, \sigma^2 = 3.2 \text{ in}^2)$ dist.
($\sigma = 3.2 \text{ in}$)

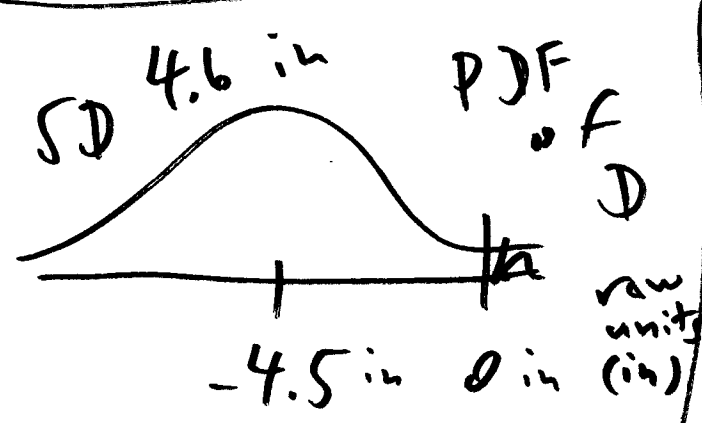
Pop. of adult u.s. men: height follows $N(\mu = 69.5 \text{ in}, \sigma^2 = 3.3 \text{ in}^2)$ dist.

1 woman chosen at random, height \underline{W} ; (265)
 1 man chosen at random (independently),
 height \underline{M} ; $P(\text{woman taller than man})$
 $= P(\underline{W} > \underline{M})$
 $= ?$

Define $D = W - M$

By consequence (4), $D \sim N(65 - 69.5 = -4.5$
 in, $3.2^2 + 3.3^2 = 21.1$
 in²)

$P(\underline{W} > \underline{M}) = P(D > 0)$



Convert to z :
 $\frac{0_{in} - (-4.5_{in})}{4.6_{in}} = +0.98$

(Z table)
 0.8365 | .1635
 (28 May 19)

So $P(\underline{W} > \underline{M}) = 16\%$
 (about 1 in 6)

Def) rv $X_1, \dots, X_n \rightarrow$ sample mean (266)

of (X_1, \dots, X_n) is $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

Consequence,
continued

$$\textcircled{5} \left\{ \begin{array}{l} X_i \stackrel{i.i.d.}{\sim} N(\mu, \sigma^2) \\ (i=1, \dots, n) \end{array} \right\}$$

$$\rightarrow \bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$\text{so } SD(\bar{X}_n) = \frac{\sigma}{\sqrt{n}}$$

Because $E(\bar{X}_n) = \mu$, \bar{X}_n is an def.

unbiased estimator of μ

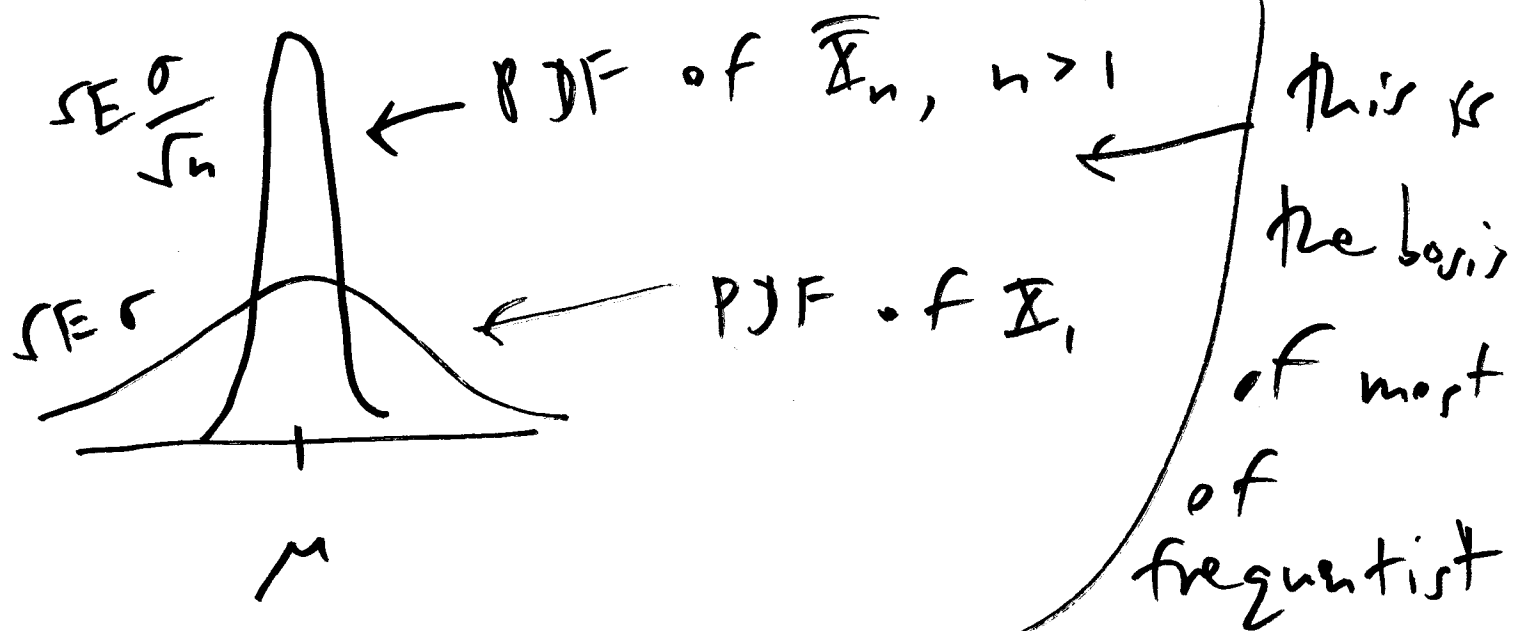
In frequentist statistics,

the standard deviation (SD) of an

estimator $\hat{\theta}_n^{(rv)}$ of a parameter θ is

called the standard error $SE(\hat{\theta})$ of $\hat{\theta}_n$

So if you use \bar{X}_n as an estimate ⁽²⁶⁷⁾ of μ , $SE(\bar{X}_n) = \frac{\sigma}{\sqrt{n}} \rightarrow 0$ as $n \rightarrow \infty$



As $n \uparrow$, \bar{X}_n gets better as an estimate of μ , at a $\frac{1}{\sqrt{n}}$ rate. This is called the square root law.

Unfortunately, this means that to cut the $SE(\bar{X}_n)$ in half (you have to quadruple the sample size).

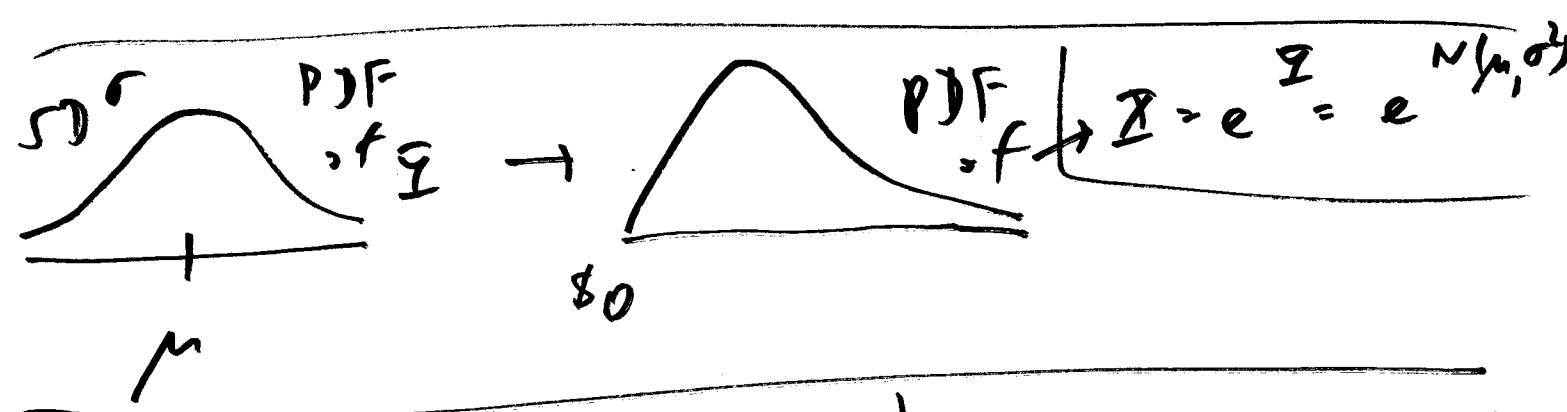
log normal Distribution } (This distribution is mis-named (1900); it should be called the

Exponential-Normal distribution, but we're stuck with a bad name.)

Def.

$X > 0$

If $Z = \log(X) \sim N(\mu, \sigma^2)$, people say that $X \sim$ Log Normal (μ, σ^2) .



$X \sim$ Log Normal (μ, σ^2)

$Z = \log(X) \sim N(\mu, \sigma^2)$

~~scribble~~
Can get MGF of X from MGF of Z

MGF of \mathbb{I} is $\psi_{\mathbb{I}}(t) = \exp(\mu t + \frac{1}{2}\sigma^2 t^2)$ (269)

But by definition

$$\psi_{\mathbb{I}}(t) = E(e^{t\mathbb{I}}) = E(e^{t \log \mathbb{X}})$$

$$= E(\mathbb{X}^t), \text{ so we can}$$

$$E(\mathbb{X}) = \psi_{\mathbb{I}}(1)$$

$$= \exp\left(\mu + \frac{\sigma^2}{2}\right)$$

read the moments of \mathbb{X} directly from the ~~MGF~~ MGF of \mathbb{I}

$$V(\mathbb{X}) = \psi_{\mathbb{I}}(2) - \left[\psi_{\mathbb{I}}(1)\right]^2$$

$$= \exp(2\mu + \sigma^2) [e^{\sigma^2} - 1]$$

Famous Case Study

~~exercise~~

(Known constant)

Pricing stock options, continued

1 share of a stock, current

price S_0 . Heroic assumption: price

u time units in the future will be e^{270}

$$S'_u = S'_0 e^{\zeta_u}, \quad \zeta_u \sim N(\mu u, \sigma^2 u).$$

Can write $S'_0 e^{\zeta_u} = e^{\zeta_u + \log(S'_0)}$. Now

$$\left[\zeta_u + \log(S'_0) \right] \sim N(\mu u + \log(S'_0), \sigma^2 u),$$

$$\text{So } S'_u \sim \text{Log Normal}[\mu u + \log(S'_0), \sigma^2 u].$$

Consider a single time horizon u ;

heroic
assumption
rewritten \rightarrow

$$S'_u = S'_0 \exp[\mu u + (\sigma\sqrt{u}) \cdot \zeta_1],$$

$$\zeta_1 \sim N(0, 1)$$

we need to price the option to buy 1 share of this stock for price q at time u .

Use risk-neutral pricing as in the (271) previous discussion: force present value

$E(S_u) \stackrel{\Delta}{=} S_0$. Let time scale of u be in years; let ^{the} risk-free (continuous-compounding) interest rate be r /year;

then present value of ~~S_u~~ S_u is $e^{-ru} \cdot E(S_u)$.

But by heroic ^{log normal} assumption, $E(S_u) = S_0 \exp(\mu u + \frac{\sigma^2 u}{2})$ so set S_0 equal to

result is $\mu = r - \frac{\sigma^2}{2}$ $e^{-ru} S_0 \exp(\mu u + \frac{\sigma^2 u}{2})$ for risk-neutral pricing.

Value of option at time u will be (272)

$h(S_u)$, where $h(S) = \begin{cases} S - g & \text{if } S > g \\ 0 & \text{else} \end{cases}$.

with $\mu = r - \frac{\sigma^2}{2}$, $h(S_u) > 0$ iff

$$Z > \frac{\log\left(\frac{g}{S_0}\right) - \left(r - \frac{\sigma^2}{2}\right)u}{\sigma\sqrt{u}} \triangleq c$$

Now a nasty integral

arise: risk-neutral price of option is the present value of $E[h(S_u)]$,

which is

$$e^{-ru} E[h(S_u)] = e^{-ru} \int_c^{\infty} \left[S_0 e^{(r - \frac{\sigma^2}{2})u + \sigma z\sqrt{u}} - g \right] \cdot$$

Careful calculation reveals the following (famous) formula:

$$\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) dz$$

$$S_0 \Phi(\sigma\sqrt{u} - c) - q e^{-ru} \Phi(-c)$$
 is

the risk-neutral price of the option,

where $c = \log\left(\frac{q}{S_0}\right) - \left(r - \frac{\sigma^2}{2}\right)u$ ← This formula

(Black-Scholes)

$\sigma\sqrt{u}$

was derived in 1973 by

Gamma Distribution

(American economist) → Fischer Black

(1938-1995) died

($\alpha, \beta > 0$) X has the

and (age 57 throat cancer)

Gamma dist. with parameters (α, β),

Canadian-American economist → Myron Scholes (1941-)

won Nobel prize

with $X \sim \Gamma(\alpha, \beta)$ or

$X \sim \text{Gamma}(\alpha, \beta)$ →

in Economics for this work in 1997, together with Robert

X continuous on $(0, \infty)$ with

American (economist) → Myron Scholes

Merton (1944-2003)

PDF $f_{\mathbb{R}}(x|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \mathbb{I}(x > 0)$

support of \mathbb{R}

α is called a shape parameter in the

$\Gamma(\alpha, \beta)$ family because it governs things like skewness of the dist.

β is related to the scale of the distribution, which measures how spread out the

dist. is $\Gamma(\alpha)$ is the Gamma function,

invented to deal with integrals of functions like \otimes above:

$$\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx$$

has no anti-derivative in closed form

(275)

$\Gamma(x)$ turns out to be a continuous generalization of the factorial function, because

$$\left(\begin{array}{c} n \text{ positive} \\ \text{integer} \end{array} \right) \rightarrow \Gamma(n) = (n-1)!$$

$\Gamma(x) \rightarrow \infty$ really quickly as $x \rightarrow \infty$, so it's better to evaluate the Gamma PDF on the log scale and then exponentiate:

$$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} = \exp \left[\alpha \ln(\beta) - \ln \Gamma(\alpha) + (\alpha-1) \ln(x) - \beta x \right]$$

Another way to tame $\Gamma(x)$ is with a Stirling's approximation:

$$\Gamma(x) \approx \sqrt{2\pi} x^{x-\frac{1}{2}} e^{-x}$$

for large x

so that $\ln f(x) = \frac{1}{2} \ln(2\pi) + (x - \frac{1}{2}) / \ln x - x$ (2)6

$X \sim I(\alpha, \beta)$ | $\psi_X(t) = \left(1 - \frac{t}{\beta}\right)^{-\alpha}$ for $t < \beta$

so $E(X) = \frac{\alpha}{\beta}$

and $V(X) = \frac{\alpha}{\beta^2}$

$SD(X) = \frac{\sqrt{\alpha}}{\beta}$

Alternative expression

$\psi_X(t) = \left(\frac{\beta}{\beta - t}\right)^{\alpha}$ for $t < \beta$

Special case of $I(\alpha, \beta)$

with $\alpha = 1$ the PDF is

$f_X(x | \beta) = \beta e^{-\beta x} I(x > 0)$

But this is just our old friend the Exponential distribution.

$X \sim \text{Exponential}(\beta)$

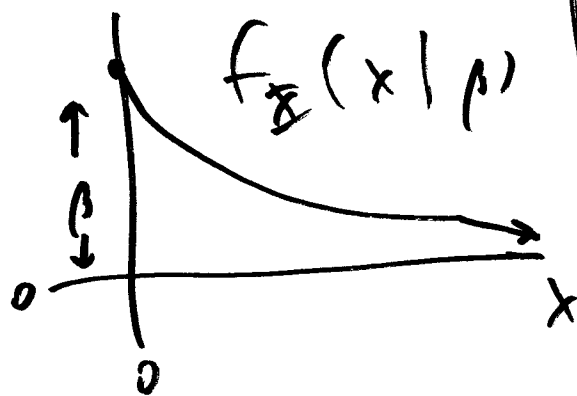
$$f_X(t) = \frac{\beta}{\beta - t}, \quad t < \beta$$

(277)

$$E(X) = \frac{1}{\beta}$$

$$V(X) = \frac{1}{\beta^2}$$

$$D(X) = \frac{1}{\beta}$$



~~Notice that the Exponential distribution has $E(X) = \frac{1}{\beta}$ equal to $D(X)$; this suggests it's related somehow to the Poisson dist.~~

Theorem Suppose

that arrivals (events) occur

according to a Poisson process with

rate β per unit time.

$$\text{and define } T_1 = T_1 - 0$$

$$T_2 = T_2 - T_1$$

$$\dots T_k = T_k - T_{k-1} \text{ for } k = 2, 3, \dots$$

Set $T_k =$ time until k^{th} arrival
 $k = 1, 2, \dots$

The T_i are called the inter-arrival (278)

times.

The Then it turns out that $T_i \stackrel{\text{IFD}}{\sim} \text{Exponential}(\beta)$

Exponential dist. is also related to the Geometric dist., in that they both

have a memoryless property Theorem

$X \sim \text{Exponential}(\beta)$; $t > 0, h > 0$

$$\rightarrow P(X \geq t+h | X \geq t) = P(X \geq h)$$

Example $X =$ ^{from initial use} time until a manufactured product fails (e.g., light bulb)

$$F_X(x) = P(X \leq x) \quad | \quad 1 - F_X(x) = P(X > x)$$

$= P(\text{"system survives" at least to time } x)$

For this reason, $1 - F_X(x)$ is called (279)

the survival function $S_X(x) = 1 - F_X(x)$

in medicine and the reliability function

$R_X(x) = 1 - F_X(x)$ in engineering.

Earlier we showed that $F_X(x) = 1 - e^{-\beta x}$
for $X \sim \text{Exponential}(\beta)$ for $x > 0$

So $S_X(x) = R_X(x) = e^{-\beta x}$ for this dist.

The instantaneous failure rate or hazard rate

function is defined to be $H_X(x) = \frac{f_X(x)}{S_X(x)}$

This gives $P(\text{failure in interval } (x, x+\epsilon) \mid \text{survival to time } x)$ for small $\epsilon > 0$ $= \frac{f_X(x)}{R_X(x)}$

Notice that if $X \sim \text{Exponential}(\beta)$ (250)

$$\text{then } H_X(x) = \frac{\beta e^{-\beta x}}{e^{-\beta x}} = \beta \left(\frac{\text{Constant in } x}{x} \right)$$

The Exponential is the only failure rate distribution with constant hazard. Returning

to the earlier result that $X \sim \text{Exponential}(\beta)$,

$$\rightarrow P(X \geq t+h | X \geq t) = P(X \geq h),$$

for all
 $t \geq 0$
 $h \geq 0$

this says that if the product has survived to time t , the chance it

will survive to time $(t+h)$ is the same as the original chance of surviving from time 0 to time h ; i.e., the

system doesn't remember how long it's survived" (this ^{often} notes the Exponential unrealistic in practice)

Consequence ① $X_i \stackrel{iid}{\sim}$ Exponential (β) (281)
($i=1, \dots, n$),

then

$Y_1 = \min(X_1, \dots, X_n) \sim \text{Exponential}(n\beta)$.

Beta $\alpha, \beta > 0$ $X \sim \text{Beta}(\alpha, \beta) \leftrightarrow$

Distribution $f_X(x) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$.

The name comes from $\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}$ support of X

the normalizing constant: the function $x^{\alpha-1} (1-x)^{\beta-1}$ has no closed-form

anti-derivative, so people just made

Definition For all $\alpha > 0$
 $\beta > 0$ $B(\alpha, \beta) \triangleq \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx$
beta function

Can show that $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$. (282)

(α, β) jointly control

the shape of the Beta (α, β) dist.

(yuck)

$X \sim \text{Beta}(\alpha, \beta)$

$$f_X(t) = 1 + \sum_{k=1}^{\infty} \left(\prod_{r=0}^{k-1} \frac{\alpha+r}{\alpha+\beta+r} \right) \frac{t^k}{k!}$$

$$E(X) = \frac{\alpha}{\alpha+\beta}$$

$$V(X) = \left(\frac{\alpha}{\alpha+\beta} \right) \left(\frac{\beta}{\alpha+\beta} \right) \left(\frac{1}{\alpha+\beta+1} \right)$$

Case Study

~~DeLoe~~

(Castaneda
v. Partida
continued)

$n=220$ grand jurors chosen from ~~(eligible)~~ eligible population of Hidalgo County, Texas, which was 79.1% Mexican-American, but only $s=100$

selected grand jurors were Mexican-American; let's summarize the information in a Bayesian fashion about evidence of discrimination.

Data $S = \#$ Mexican-American ^{chosen} in jury selection of $n = 220$ people 283

Unknown $\theta =$ actual probability of an eligible Mexican-American person being chosen ($0 < \theta < 1$)

Sampling Model $(S | \theta) \sim \text{Binomial}(n, \theta)$,

\leftarrow PMF

i.e., $f_{S|\theta}(s|\theta) = P(S=s|\theta) = \binom{n}{s} \theta^s (1-\theta)^{n-s}$. $I(s=0, 1, \dots, n)$

Bayesian approach ① Information internal to dataset about θ summarized

by the likelihood (un-normalized) density, defined to be $l(\theta | s) = c P(S=s|\theta)$,

c an arbitrary positive constant — ^{just} think of $P(S=s|\theta)$ as a function of θ for ^{fixed} s

Here $l(\theta | s) = c \binom{h}{s} \theta^s (1-\theta)^{h-s}$ can be absorbed into c since c does not depend on θ

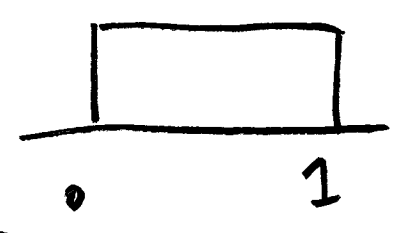
$= c \theta^s (1-\theta)^{h-s}$

(2) Information external to dataset about θ summarized by the prior density $f(\theta)$.

Here are some

possibilities for the prior, depending on your knowledge base:

(a) neutral prior $\theta \sim \text{Uniform}(0,1)$



this dist. embodies the information { θ could be anywhere between 0 and 1, with no value favored }

(b) cut the district attorney some slack prior



this prior gives the DA the benefit of the doubt

When you're uncertain about what prior 285 to use, write down all the reasonable priors & do a sensitivity analysis (use each prior one by one & see if ^{posterior} answer is the same) essentially

③ Combine internal & external information

with
Bayes'
Theorem

$$f_{\Theta|S}(\theta|s) = c \cdot f_{\Theta}(\theta) \cdot f(\theta|s)$$

↑
↑
↑
↑

posterior (information)
=
(normalizing constant)
·
(prior information)

·
(likelihood information)

Here

$$f_{\Theta|S}(\theta|s) = c f_{\Theta}(\theta) \theta^s (1-\theta)^{n-s}$$

Rev. Bayes himself noticed back in 1760

that if you take $f_{\theta}(\theta) = c \theta^{\text{power}_1} (1-\theta)^{\text{power}_2}$
 then the product of 2 such densities is
 another such density, meaning that the
 posterior would have the same form as
 the prior & likelihood, making calculations

easier

Moreover, we already know the
 name of densities that look like $\theta^{\text{power}_1} (1-\theta)^{\text{power}_2}$:

The $X \sim \text{Beta}(\alpha, \beta)$ ($\alpha > 0, \beta > 0$) \rightarrow

Beta
 distributions $f_X(x) = c \theta^{\alpha-1} (1-\theta)^{\beta-1}$

as our prior PDF

So let's take $f_{\theta}(\theta) = c \theta^{\alpha-1} (1-\theta)^{\beta-1}$

in the law suit case study; then

$$f_{\theta|S}(\theta|s) = c \left[\theta^{\alpha-1} (1-\theta)^{\beta-1} \right] \left[\theta^s (1-\theta)^{n-s} \right]$$

$$= c \theta^{(\alpha+s)-1} (1-\theta)^{(\beta+n-s)-1} = \text{Beta}(\alpha+s, \beta+n-s)$$

So the prior-to-posterior updating looks like this:

Beta dist. is conjugate to the Binomial likelihood

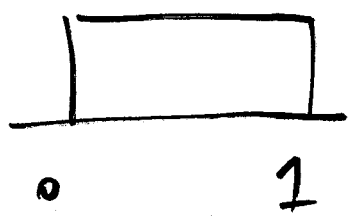
$$\left. \begin{aligned} \theta &\sim \text{Beta}(\alpha, \beta) \\ (s' | \theta) &\sim \text{Binomial}(n, \theta) \end{aligned} \right\} \rightarrow (\theta | s) \sim \text{Beta}(\alpha+s, \beta+n-s)$$

$s = 100$
 $n = 220$

How choose (α, β) ?

(a) Neutral prior

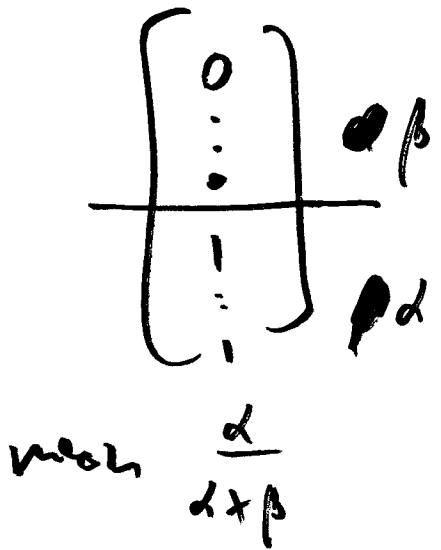
but $\text{Uniform}(0, 1) = \theta^{1-1} (1-\theta)^{1-1}$



So $\theta \sim \text{Uniform}(0, 1) \leftrightarrow \theta \sim \text{Beta}(1, 1)$

(b) cut DA stack prior

There's an extremely useful thing that happens with conjugate priors:

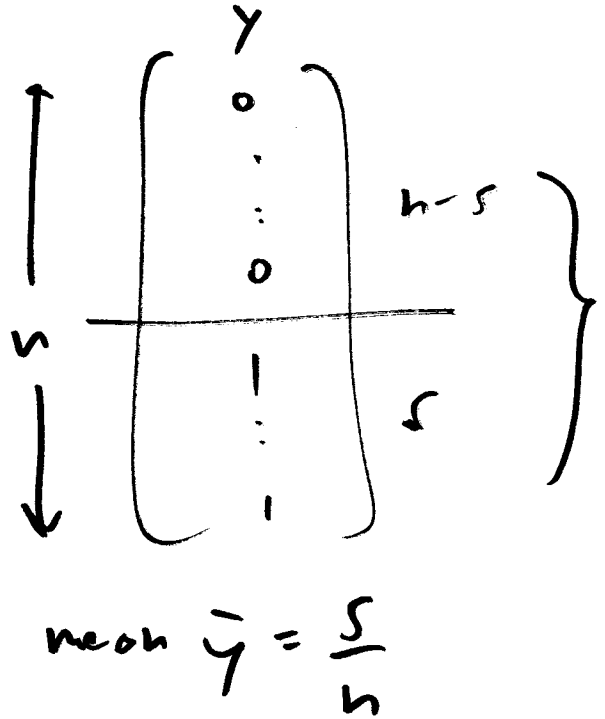


pseudo data

prior effective sample size $(\alpha + \beta)$

Beta prior distribution acts like a dataset with α 1s & β 0s

with the property that



sample data

dataset sample size n

if you do a Bayesian analysis with the Beta (α, β) prior and I do a frequentist

analysis on the dataset with $(\alpha + s)$ 1s and $(\beta + h - s)$ 0s formed by merging the prior & sample datasets, we'll get the same results.

(b) Cut the JA stock prior

mean of Beta(α, β) dist. is $\frac{\alpha}{\alpha + \beta}$ (289)

$\frac{\alpha}{\alpha + \beta}$; set this equal to 0.791

Suppose I want to put in ^{prior} information equivalent to a prior sample size $\frac{1}{10}$ as big as the data sample size (507); set

$$(\alpha + \beta) = \frac{1}{10} n = 22$$

$$\text{solve: } \begin{cases} \alpha = 17.4 \\ \beta = 4.6 \end{cases}$$

$$n = 220$$

$$s = 100$$

likelihood is

$$c \theta^s (1-\theta)^{n-s} = c \theta^{(s+1)-1} (1-\theta)^{(n-s+1)-1}$$

= Beta($s+1, n-s+1$) dist

(101)

(121)

(a) Neutral prior:

$$\text{Beta}(1, 1)$$

prior sample size 2

prior is

$$\text{Beta}(\alpha + s, \beta + n - s)$$

\uparrow
101

\uparrow
(121)

(same as likelihood)

(b) cut
DA
stock
prior

Beta (17.4, 4.6) prior
 α β

(290)

posterior \rightarrow Beta ($\alpha + s$, $\beta + n - s$)
 \uparrow \uparrow
 (117.4) (124.6)

220 100
 \downarrow \downarrow

prior	posterior		posterior mean of θ is $\frac{\alpha + s}{\alpha + \beta + n}$
	mean	SD	
neutral	0.455	0.0333	
cut DA stock	0.485	0.0321	

Posterior SD is $\sqrt{\left(\frac{\alpha + s}{\alpha + \beta + n}\right) \left(\frac{\beta + n - s}{\alpha + \beta + n}\right) \left(\frac{1}{\alpha + \beta + n + 1}\right)}$

The no-discrimination rate of 0.791 is

$\frac{0.791 - 0.455}{0.0333} = 10.1$ posterior SDs away from posterior expectation

under the neutral prior and

$$\frac{0.791 - 0.485}{0.0321} = 9.5 \text{ posterior S.D.s}$$

away from posterior expectation under the cut-DA-slack prior; there was Q.E.D. discrimination

Multinomial Distributions (back to discrete) / You're contemplating a population that contains elements of $k \geq 2$ types (e.g., {Democrat, Republican, Libertarian, Independent, Green, ^{other}}).

Suppose the proportion of elements of type i is $0 \leq p_i \leq 1$ with $\sum_{i=1}^k p_i = 1$; $\mathbf{p} = (p_1, \dots, p_k)$.

You take an IID sample of size n (292)
 from this pop.; $X_i = \#$ elements of
 type i in your sample; $\sum_{i=1}^k X_i = n$.

Can show that the vector $\underline{X} = (X_1, \dots, X_k)$

has
 M
 $P.F.$

$$f_{\underline{X}|n, \underline{p}}(x_1, \dots, x_k) = \begin{cases} \binom{n}{x_1, \dots, x_k} p_1^{x_1} \dots p_k^{x_k} & \text{if } \sum_{i=1}^k x_i = n \\ 0 & \text{else} \end{cases}$$

where $\left(\sum_{i=1}^k p_i = 1 \right)$

$$\binom{n}{x_1, \dots, x_k} \triangleq \frac{n!}{x_1! x_2! \dots x_k!}$$

is the multinomial coefficient

This is called the Multinomial (n, \underline{p})
 distribution.

$$E(X_i) = np_i \quad V(X_i) = np_i(1-p_i)$$

(just like binomial)

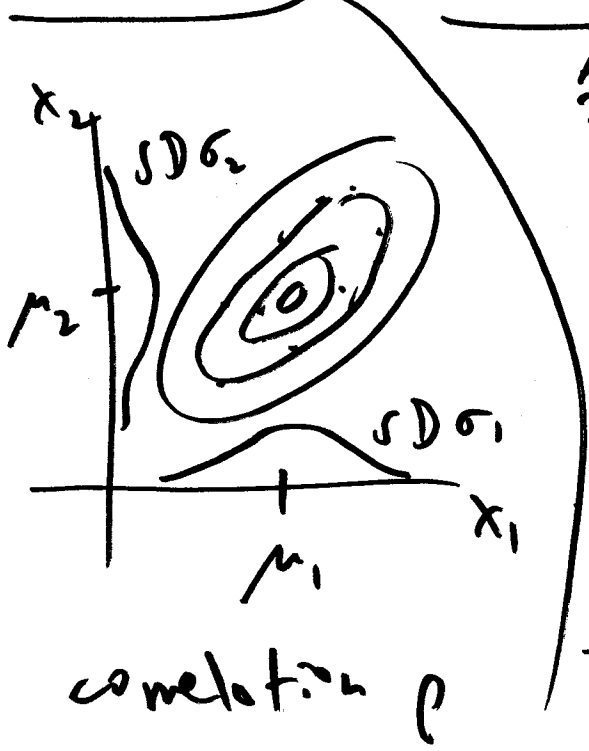
But now something new:

$$C(X_i, X_j) = -n p_i p_j$$

negatively correlated because $\sum_{i=1}^k X_i = n$

Bivariate Normal Dist.

Can build a 2-dimensional (bivariate) version of the Normal dist. as follows:



$$Z_1, Z_2 \stackrel{IID}{\sim} N(0, 1)$$

Specify 5 parameters:

$-\infty < \mu_1 < +\infty$	$0 < \sigma_1 < \infty$
$-\infty < \mu_2 < +\infty$	$0 < \sigma_2 < \infty$
$-1 < \rho < +1$	

Now build (X_1, X_2) with the transformation $X_1 = \mu_1 + \sigma_1 Z_1$

$$X_2 = \sigma_2 \left[\rho Z_1 + \sqrt{1-\rho^2} Z_2 \right] + \mu_2$$

The joint PDF of $\underline{X} = (X_1, X_2)$ is

then $f_{X_1, X_2}(x_1, x_2) = \frac{1}{2\pi\sqrt{1-\rho^2}\sigma_1\sigma_2} \cdot \exp \left\{ \right.$

$$-\frac{1}{2(1-\rho^2)} \left[\left(\frac{x_1 - \mu_1}{\sigma_1} \right)^2 - 2\rho \left(\frac{x_1 - \mu_1}{\sigma_1} \right) \left(\frac{x_2 - \mu_2}{\sigma_2} \right) \right.$$

standard units $\left. + \left(\frac{x_2 - \mu_2}{\sigma_2} \right)^2 \right]$

This is the Bivariate normal $(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$ dist.

Easy to show that $E(X_1) = \mu_1$, (295)

$$E(X_2) = \mu_2, \quad V(X_1) = \sigma_1^2, \quad V(X_2) = \sigma_2^2,$$

$$\rho(X_1, X_2) = \rho.$$

Consequences of this def.

① $(X_1, X_2) \sim \text{Bivariate Normal} \rightarrow$

$$\left(\begin{array}{l} X_1, X_2 \\ \text{independent} \end{array} \right) \leftrightarrow \left(\begin{array}{l} X_1, X_2 \\ \text{uncorrelated} \end{array} \right)$$

we already knew the \rightarrow direction is general; what's new here is that correlation 0 implies independence

if $(X_1, X_2) \sim \text{Bivariate Normal}$.

② $(X_1, X_2) \sim \text{Bivariate Normal}(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$

→ conditional distribution of X_2

given that $X_1 = x_1$ is (univariate)

normal with mean $E(X_2 | x_1) =$

$$\mu_2 + \frac{\rho \sigma_2}{\sigma_1} (x_1 - \mu_1)$$

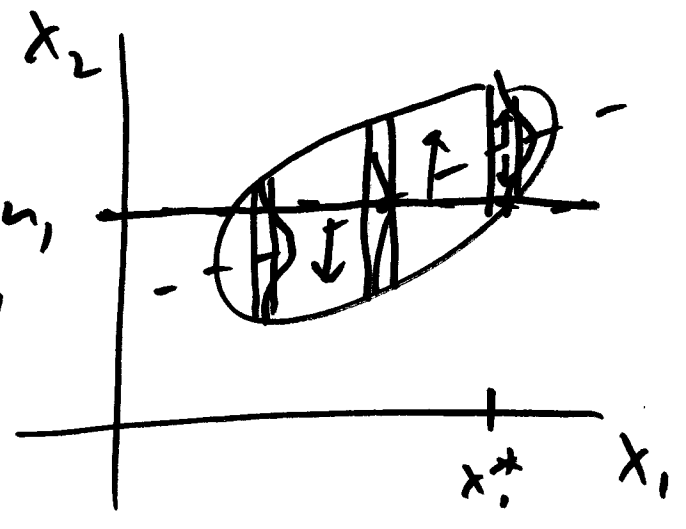
and variance $V(X_2 | x_1)$

$$= (1 - \rho^2) \sigma_2^2$$

above

Result ② says that if (X_1, X_2) are

Galton, verified



conditional

Bivariate Normal then the distributions of X_2 given $X_1 = x_1^*$ in all of the vertical strips are also normal.

And the means of all these normal distributions in the vertical strips are connected together by Galton's

regression
line

$$\hat{x}_2 = \mu_2 + \left(\frac{\rho\sigma_2}{\sigma_1}\right)(x_1 - \mu_1)$$

This line has slope $\beta_1 = \frac{\rho\sigma_2}{\sigma_1}$ and "y"-intercept

$$\beta_0 = \mu_2 - \beta_1\mu_1$$

Moreover,

$$\hat{x}_2 = \beta_0 + \beta_1 x_1$$

we can now quantify an earlier insight:

ignore x_1 ,

$$\text{predict } (\hat{x}_2)_{x_1} = \mu_2 = E(X_2)$$

(root mean squared error)

(RMSE) of this prediction is

$$\sqrt{V(X_2)} = \sigma_2$$

Use x_1
to predict
 x_2

$$\text{pred. 2} + \left(\hat{x}_2\right)_{\text{use } x_1} = E(X_2 | X_1 = x_1)$$

$$= \mu_2 + \frac{\rho \sigma_2}{\sigma_1} (x_1 - \mu_1)$$

RMSE of this

prediction is $\sqrt{V(X_2 | x_1)} = \sigma_2 \sqrt{1 - \rho^2}$

~~RMSE(used) \leq RMSE(ignoring x_1)~~

Since $-1 < \rho < 1$, $\sigma_2 \sqrt{1 - \rho^2} \leq \sigma_2$

with equality only when $\rho = 0$.

③ $(X_1, X_2) \sim \text{Bivariate Normal}(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$

$$Y = a_1 X_1 + a_2 X_2 + b, \quad (a_1, a_2, b) \text{ arbitrary constants}$$

$$\rightarrow Y \sim N(a_1 \mu_1 + a_2 \mu_2 + b, a_1^2 \sigma_1^2 + a_2^2 \sigma_2^2 + 2a_1 a_2 \rho \sigma_1 \sigma_2)$$

Large
Random
Samples

(DS ch. 6)

You draw an IID random sample X_1, \dots, X_n from a population with the goal of estimating the population mean $\mu = E(X_i)$.

We've already seen that, from a root mean squared error point of view, the sample mean $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ is the best you can do (in the absence of prior information).

It would be nice if \bar{X}_n approached the

right answer μ as n increases; how to quantify that idea?

Two inequalities that help

Markov inequality

(300)

Suppose

X is a non-negative r.v., i.e.

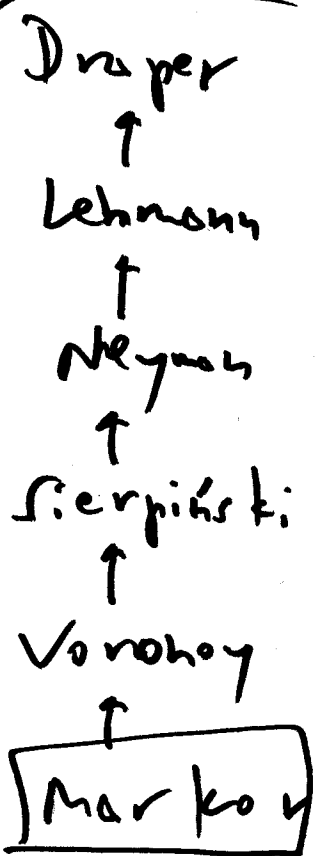
$$P(X \geq 0) = 1$$

then for all

$$\text{real } t > 0, \quad P(X \geq t) \leq \frac{E(X)}{t} \quad *$$

(Attributed to Andrey Markov (1856-1922), a Russian mathematician who did pioneering work on stochastic processes)

* Says that, if $E(X)$ is fixed, you can't move more & more probability out into the right tail beyond a certain point.



25 April

Laplace

