

It turns out that e^{-cx^2} has no 260
(c>0) 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, $\Phi(-x) = 1 - \Phi(x)$

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

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

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

and $F_{\frac{X-\mu}{\sigma}}^{-1}(p) = \mu + \sigma \Phi^{-1}(p)$

Empirical Rule

Part 1 start at the mean μ 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 Diff. 2SDs

either way: $(\mu \pm 2\sigma)$ captures (about

95%) of the probability

Part 3

Diff. 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:

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

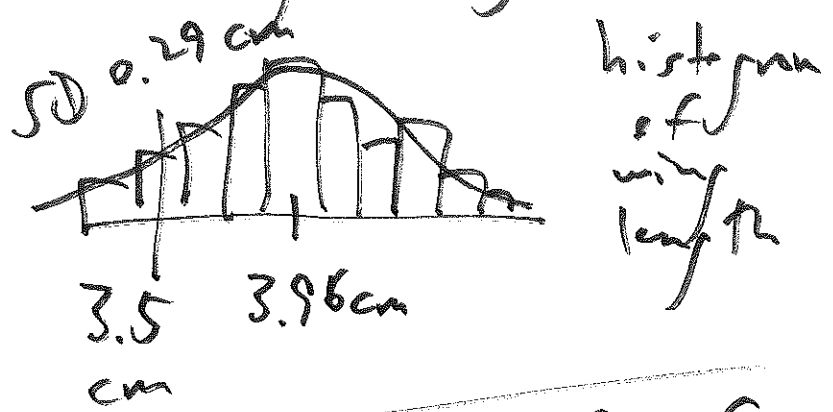
$y_1 = 4.1$

$y_2 = 3.3$

\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

3.2
3.3
:
3.5
3.5
3.5
3.5
3.6
:
:
4.7

↑

8

↓

↑

103

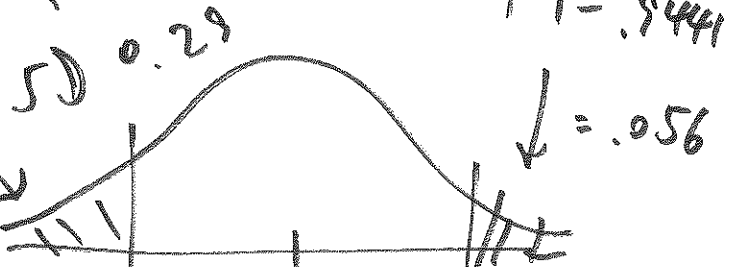
↓

A_1

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

A_2 (approximate)

← .9441 → 1 - .9441



3.5 3.96
mean

raw units

standard units
-1.59 0 +1.59

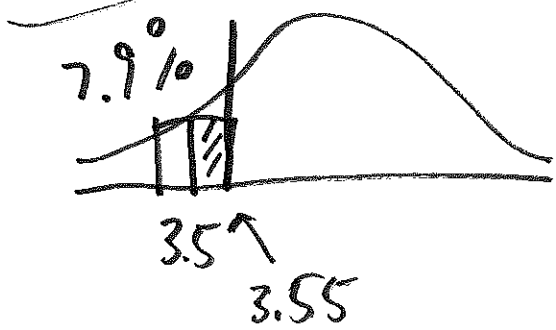
$\frac{3.5 - 3.96}{0.29}$

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, (264)
 $X_i \sim \text{Normal}(\mu_i, \sigma_i^2)$

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

nice additive property

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

Notation

$$\text{Normal}(\mu, \sigma^2) \stackrel{\Delta}{=} N(\mu, \sigma^2)$$

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

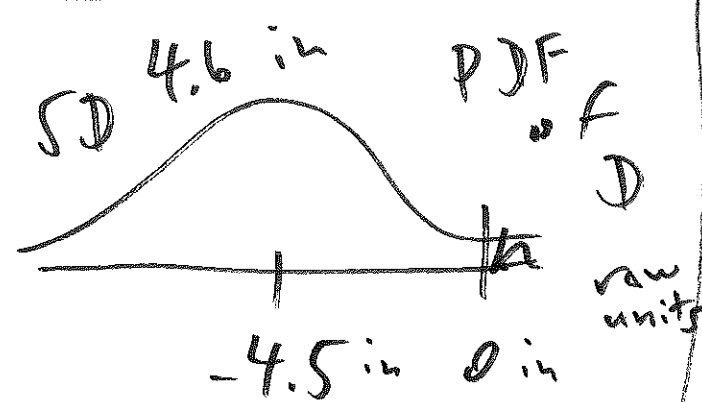
$$N(\mu = 69.5 \text{ in}, \sigma^2 = 3.3 \text{ in}^2) \text{ dist.}$$

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 $D = W - M$

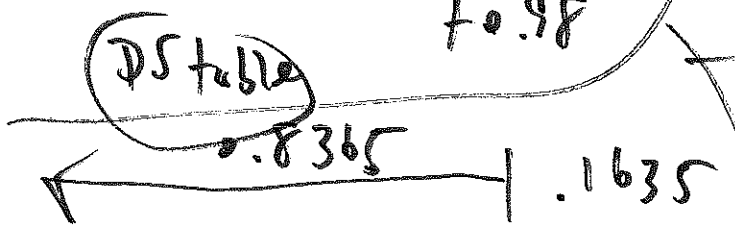
By consequence (4), $D \sim N(65 - 69.5 = -4.5$
 in, $3.2^2 + 3.3^2 = 21.1$
 in²)

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



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

SD $\sqrt{21.1 \text{ in}^2} \approx 4.6 \text{ in}$



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

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

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

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

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

unbiased estimator of μ

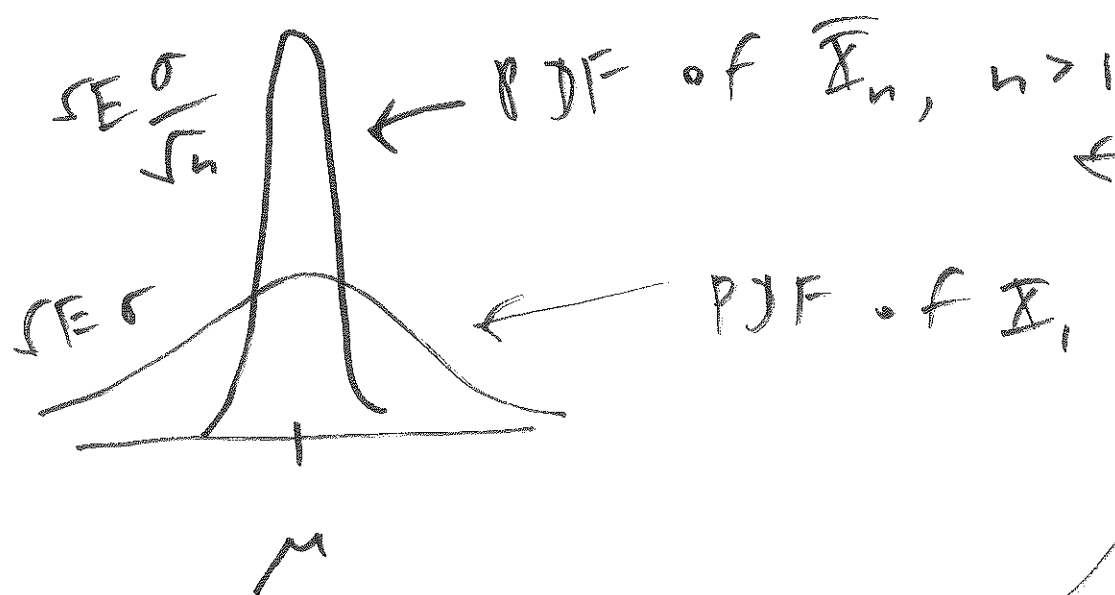
In frequentist statistics,

the standard deviation (SD) of an

estimator $\hat{\theta}_n$ 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 ^(2.67) of μ , $SE(\bar{X}_n) = \frac{\sigma}{\sqrt{n}} \rightarrow 0$ as $n \rightarrow \infty$



This is the basis of most of frequentist statistical inference

As $n \uparrow$, \bar{X}_n gets better

as an estimate of μ , at a \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.

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

But by definition

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

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

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

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

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

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

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

Famous Case Study

~~example~~

(Known constant)

price S_0 . Heroic assumption: price

Pricing stock options, continued

1 share of a stock, current

u time units in the future will be 270

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

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

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

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

Consider a single time horizon u ;

heroic
assumption
reworded \rightarrow

$$S'_u = S'_0 \exp[\mu u + (\sigma\sqrt{u}) \cdot \xi_1],$$

$$\xi_1 \sim N(0, 1)$$

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 risk-free (continuous-compounding) interest rate be (r) /year;

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

But by log normal heroic assumption,

$$E(S_u) = S_0 \exp\left(\mu u + \frac{\sigma^2 u}{2}\right)$$

S_0 set
 S_0 equal
to

result is $\mu = r - \frac{\sigma^2}{2}$ $e^{-ru} S_0 \exp\left(\mu u + \frac{\sigma^2 u}{2}\right)$
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

$$\frac{1}{2} > \frac{\log\left(\frac{g}{S_0}\right) - \left(r - \frac{\sigma^2}{2}\right)u}{\sigma\sqrt{u}} = c$$

Now a nasty integral

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

which

is

$$e^{-ru} E[h(S_u)] = e^{-ru} \int_c^{\infty} \left[S_0 e^{(r - \frac{\sigma^2}{2})u + \sigma z \sqrt{u}} - g \right] \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) dz$$

Careful calculation

reveals the (famous) formula

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

$$S_0 \mathbb{I}(\sigma\sqrt{u} - c) - q e^{-ru} \mathbb{I}(-c)$$
 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 formula)

was derived in 1973 by

Gamma Distribution

(American economist)

Fischer Black (1938-1995)

$(\alpha, \beta > 0)$ \mathbb{I} has the

Gamma dist. with parameters (α, β) ,

Canadian-American economist

Myron Scholes (1941-)

with $\mathbb{I} \sim \Gamma(\alpha, \beta)$ or

$\mathbb{I} \sim \text{Gamma}(\alpha, \beta) \rightarrow$

won Nobel prize

in Economics for this work

in 1997, together with Robert

\mathbb{I} continuous on $(0, \infty)$ with

American economist

Merton (1944-2003)

PDF $f_X(x|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} I(x > 0)$ (274)

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 (*) above:

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

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 $\ln \Gamma(x) = \frac{1}{2} \ln(2\pi) + (x - \frac{1}{2}) \ln x - x$ (276)

$X \sim \Gamma(\alpha, \beta)$ $\psi_X(t) = \left(1 - \frac{t}{\beta}\right)^{-\alpha}$ 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 $\Gamma(\alpha, \beta)$ 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.

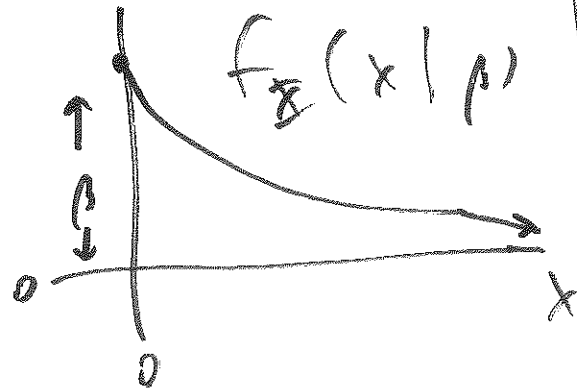
$X \sim \text{Exponential}(\beta)$

$\psi_X(t) = \frac{\beta}{\beta - t}, t < \beta$ (277)

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

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

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



~~Notice that the Exponential distribution has TMR equal to 1; 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.

and define $T_1 = T_1 - 0$

$T_2 = T_2 - T_1$

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

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

The T_i are called the inter-arrival (278)

times.

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

The

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 (eg., 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 ϵ $= \frac{f_X(x)}{R_X(x)}$

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

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

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} makes the Exponential unrealistic in practice)

Consequence ① $X_i \stackrel{i.i.d.}{\sim}$ 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 $\alpha > 0$
 $\beta > 0$ $B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx$
 \uparrow
beta function

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

(α, β) jointly control

$$\Gamma(\alpha+\beta)$$

the shape of the Beta(α, β) dist.

(yuck)

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

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

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

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

Case Study

~~Dist~~

(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;
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 $(S' | \theta) \sim \text{Binomial}(n, \theta)$,

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

Bayesian approach $\textcircled{1}$ Information internal to data set 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 $l(\theta | s) = c \binom{4}{s} \theta^s (1-\theta)^{4-s}$ can be absorbed into c since c does not depend on θ

$$= c \theta^s (1-\theta)^{4-s}$$

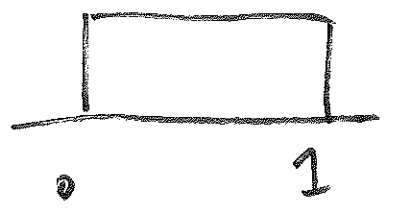
(2) Information external to dataset about θ summarized by the prior density $f(\theta)$.

Here are some

possibilities for the prior, depending on (information)

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)

③ Combine internal & external information

with
 Bayes'
 Theorem

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

\uparrow \uparrow \uparrow
 posterior (information) = (normalizing constant) · (prior information) · (likelihood information)

Here

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

Rev. Bayes himself noticed back in 1760

that if you take $f_{\theta}(\theta) = c \theta^{\text{power}} (1-\theta)^{\text{power}}$ 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 $\theta^{\text{power}} (1-\theta)^{\text{power}}$:

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

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

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

in the lawsuit case study; then

$$f_{\theta|s}(\theta|s) = c \left[\theta^{\alpha-1} (1-\theta)^{\beta-1} \right] \left[\theta^s (1-\theta)^{4-s} \right]$$

$$= c \theta^{(\alpha+s)-1} (1-\theta)^{(\beta+n-s)-1} = \text{Beta}(\alpha+s, \beta+n-s) \quad (287)$$

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

Beta dist. is conjugate to the Binomial likelihood

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

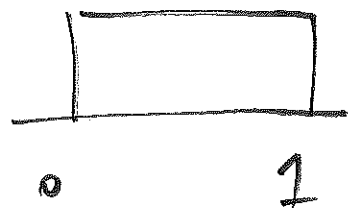
$$s = 100$$

$$n = 220$$

How choose (α, β) ?

(a) Neutral prior

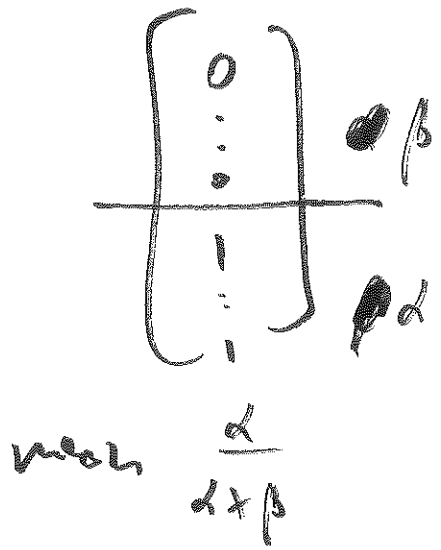
$$\text{but Uniform}(0, 1) = \theta^{1-1} (1-\theta)^{1-1}$$



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

(b) cut
DA
stock
prior

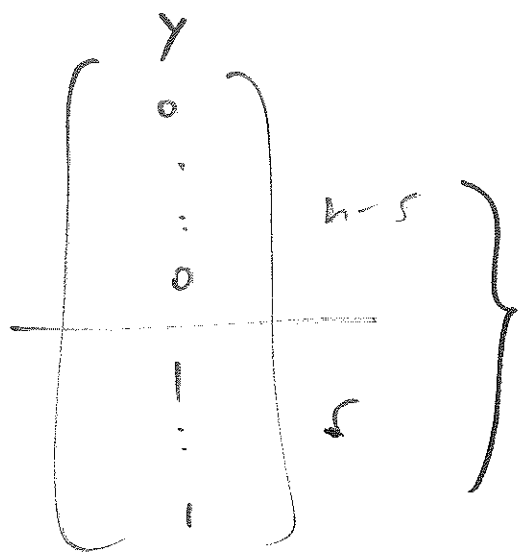
There's an extremely useful thing that happens with conjugate priors:



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

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

with the property that



data set sample size h

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

mean $\bar{y} = \frac{s}{h}$

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

(b) cut the JA stock prior

mean of Beta(α, β) dist. is $\frac{\alpha}{\alpha + \beta}$; set this equal to 0.791

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

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

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

$n = 220$
 $s = 100$

likelihood is

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

$$= \text{Beta}(s+1, n-s+1) \text{ dist}$$

(a) Neutral prior: Beta(1,1)

posterior

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

prior sample size 2

(same as likelihood)

(b) cut
DA
stock
prior

Beta (17.4, 4.6) prior
 α β

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

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

(291)

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

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}).

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

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

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

$$\binom{n}{x_1, \dots, x_k} \triangleq \frac{n!}{x_1! x_2! \dots x_k!} \text{ is the multinomial coefficient}$$

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

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

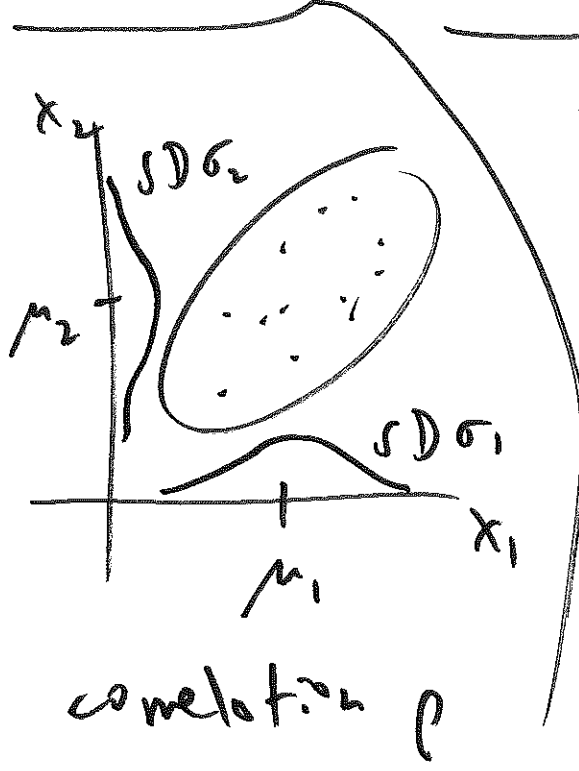
(just like binomial) But now something new:

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

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 $X_1 = \mu_1 + \sigma_1 Z_1$

(294)

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

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

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

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

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)

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

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

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 $(X_1, X_2) \sim \text{Bivariate Normal}(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$ (296)

→ 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

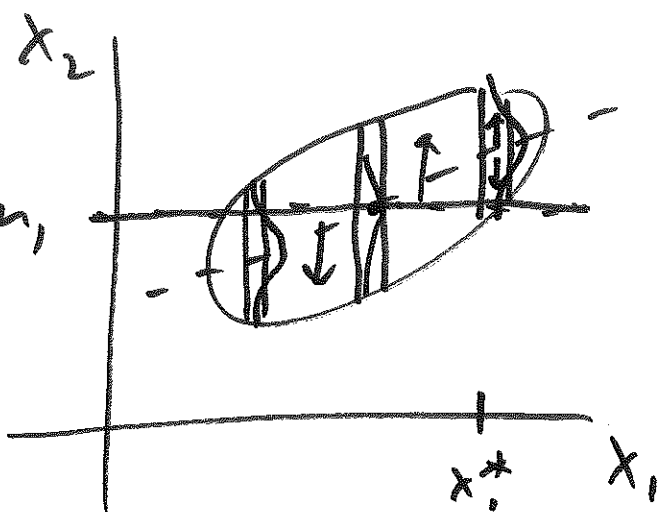
Result ② says

that if (X_1, X_2) are

Bivariate Normal then the distributions of X_2 given $X_1 = x_1^*$ in all of the

vertical strips are also normal

Galton,
visited



conditional

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

regression line

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

This line has slope $\beta_1 = \frac{\sigma_2}{\sigma_1}$ 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

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

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

RMSE of this

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

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

with equality only when $\rho = 0$.

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

$Y = a_1 X_1 + a_2 X_2 + b$, (a_1, a_2, b) 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 worst case 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

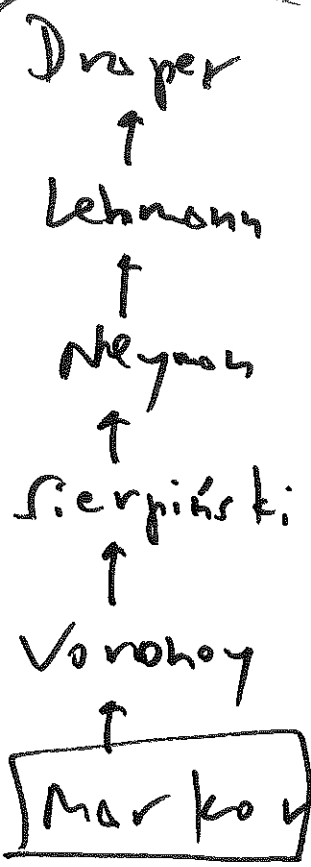
Suppose X is a non-negative r.v., i.e. $P(X \geq 0) = 1$

300

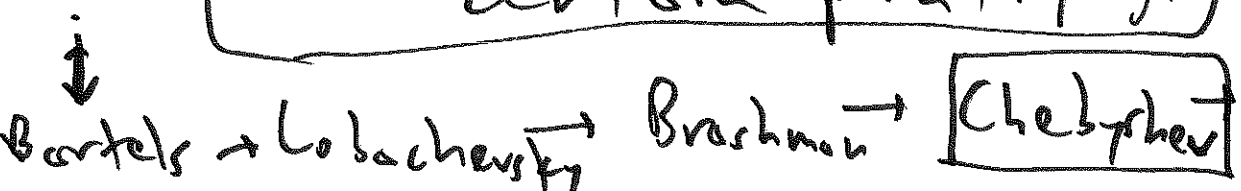
Then for all real $t > 0$, $P(X \geq t) \leq \frac{E(X)}{t}$ *

(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.



Laplace



25 April