

so that  $\ln f(x) = \frac{1}{2} \ln(2\pi) + (x - \frac{1}{2}) / \ln x - x$  (276)

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

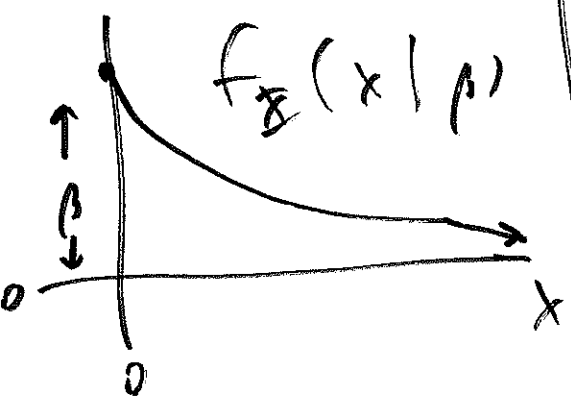
$X \sim \text{Exponential}(\beta)$

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

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

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

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



~~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  $\beta$  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  $k^{\text{th}}$  arrival  
 $k = 1, 2, \dots$

The  $T_i$  are called the inter-arrival (278)

times.

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

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

have a memory less 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)  $X =$  time <sup>from initial use</sup> until a manufactured product fails (eg., light bulb)

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

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

For this reason,  $1 - F_X(x)$  is called  
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  $\epsilon$   $= \frac{f_X(x)}{R_X(x)}$

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

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

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 <sup>often</sup> makes the Exponential unrealistic in practice)

Consequence ①  $X_i \stackrel{i.i.d.}{\sim}$  Exponential( $\beta$ ) (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

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

$\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}$   
support of  $X$

The name comes from

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)

$(\alpha, \beta)$  jointly control

the shape of the Beta  $(\alpha, \beta)$  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 <sup>chosen</sup> 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

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

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

Information internal to data set about  $\theta$  summarized

by the likelihood (un-normalized) density,

defined to be  $l(\theta | s) = c P(S=s | \theta)$

$c$  an arbitrary positive constant - <sup>just</sup> think of  $P(S=s | \theta)$  as a function of  $\theta$  for <sup>fixed</sup>  $s$ .



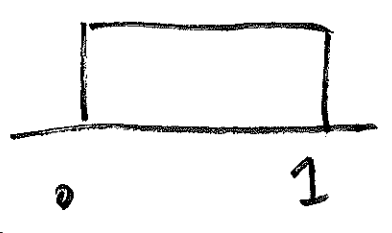
Here  $l(\theta | s) = c \binom{h}{s} \theta^s (1-\theta)^{h-s}$  could be absorbed into  $c$  since they do not depend on  $\theta$

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

(2) Information external to dataset about  $\theta$  summarized by the prior density  $f_{\theta}(\theta)$ . Here are some

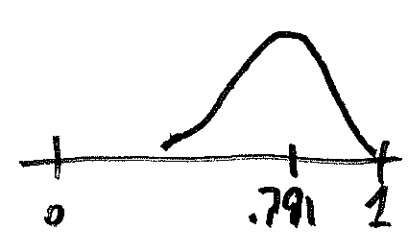
possibilities for the prior, depending on your knowledge base:

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



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

(b) but the district attorney somewhat 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 <sup>posterior</sup> 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

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

Rev. Bayes himself entitled back in 1760

that if you take  $f_{\theta}(x) = c \theta^{\text{power}} (1-\theta)^{\text{power}}$  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  $f_X(x) = c \theta^{\alpha-1} (1-\theta)^{\beta-1}$

as our prior PDF

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

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

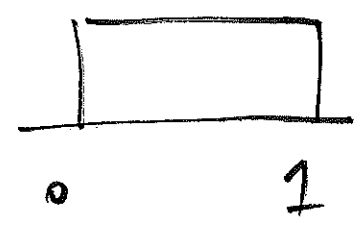
$$\left. \begin{aligned} \theta &\sim \text{Beta}(\alpha, \beta) \\ (S' | \theta) &\sim \text{Binomial}(n, \theta) \end{aligned} \right\} \rightarrow (\theta | S) \sim \text{Beta}(\alpha+s, \beta+n-s)$$

$s = 100$   
 $n = 220$

How choose  $(\alpha, \beta)$ ?

(a) Neutral prior

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



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

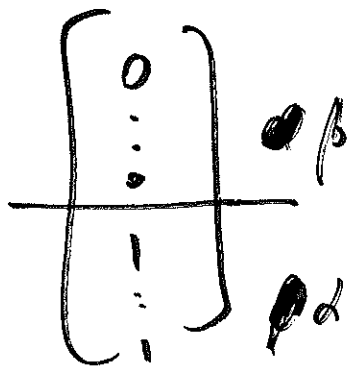
(b) cut DA stock prior

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

Beta prior distribution acts like a dataset with  $\alpha$  1s &  $\beta$  0s

pseudo data

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



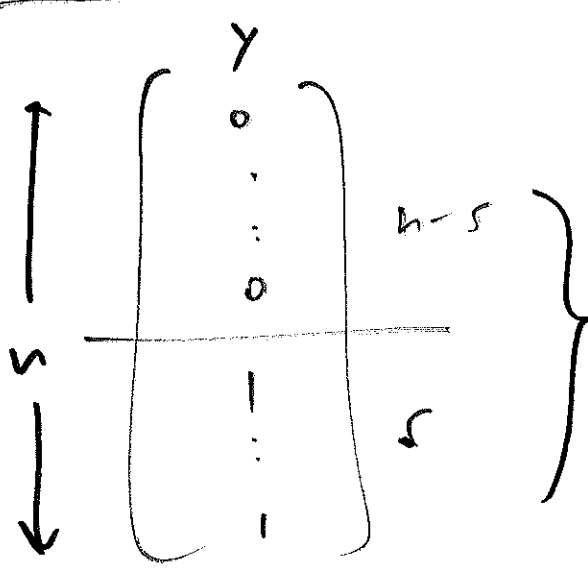
mean  $\frac{\alpha}{\alpha + \beta}$

with the property that

if you do a Bayesian analysis with the Beta( $\alpha, \beta$ ) prior and I do a frequentist

sample data

dataset sample size  $n$



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

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

(b) Cut the JA stock prior

mean of Beta( $\alpha, \beta$ ) dist. is  $\frac{\alpha}{\alpha + \beta}$  (2f9)

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

Suppose I want to put in <sup>prior</sup> 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

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

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

(101) (121)

prior

$$\text{is Beta} \left( \underset{\uparrow}{\alpha+s}, \underset{\uparrow}{\beta+n-s} \right)$$

101 (121)

(same as likelihood)

(a) Neutral prior:

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

prior sample size 2

(b) cut  
DA  
stock  
prior

Beta ( $\alpha, \beta$ ) prior  
 $\alpha$        $\beta$

posterior  
→

Beta ( $\alpha + s, \beta + n - s$ )  
↑  
 $\alpha = 117.4$

$\beta = 220$   
↓  
 $\beta + n - s = 124.6$

prior	posterior	
	mean	SD
neutral	0.455	0.0333
cut DA stock	0.485	0.0321

posterior mean of  $\theta$  is  $\frac{\alpha + s}{\alpha + \beta + n}$

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

$$\frac{0.791 - 0.485}{0.0321} = 9.5 \text{ posterior SDR}$$

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, <sup>other</sup>}).

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

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

This is called the Multinomial  $(n, p)$   
 distribution.

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

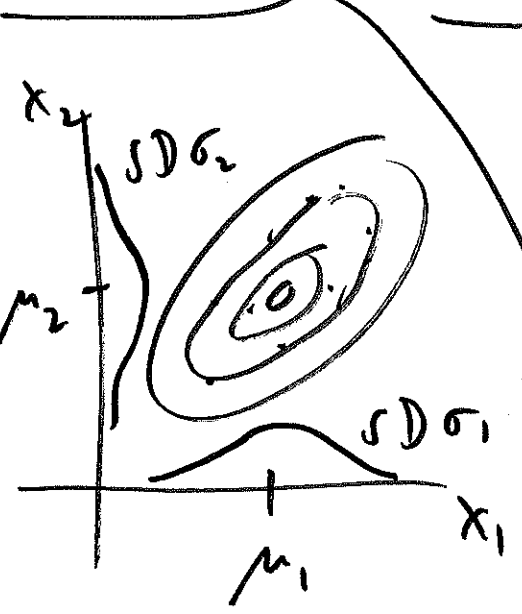
(just like binomial) But now something new:

$$C(X_i, X_j) = -np_i p_j$$

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:

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

correlation  $\rho$

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

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

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

standard units

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 in general; what's new here is that correlation 0 implies independence

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

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

$\rightarrow$  conditional distribution of  $X_2$

given that  $X_1 = x_1$  is (univariate)

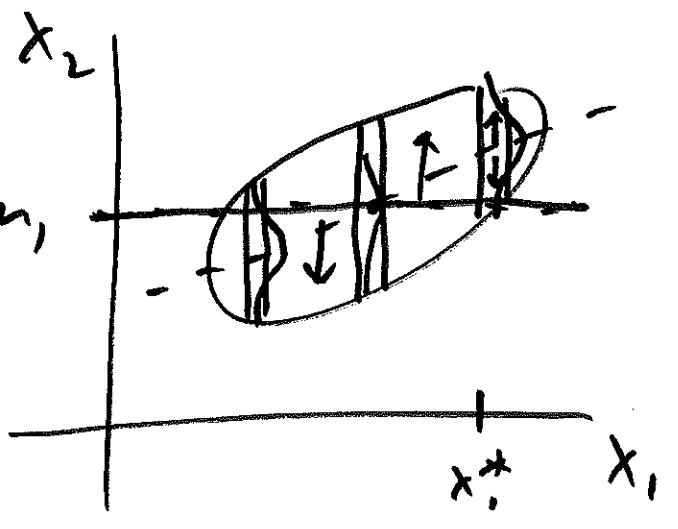
normal with mean  $E(X_2 | x_1) =$

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

and variance  $V(X_2 | x_1)$

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

Galton, revisited



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

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

regression line

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

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  $x_1$ ,

$$\text{predict } (\hat{x}_2)_{no\ x_1} = \mu_2 = E(X_2)$$

(root mean squared error)

(RMSE) of this prediction is

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

use  $x_1$   
to predict  
 $x_2$

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

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

RMSE of this

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

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  $\mu$  as  $n$  increases; how to quantify that idea?



Two inequalities that help

# Markov inequality

Suppose

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

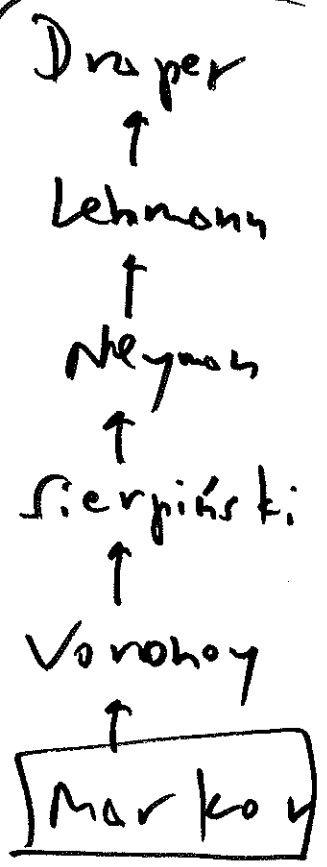
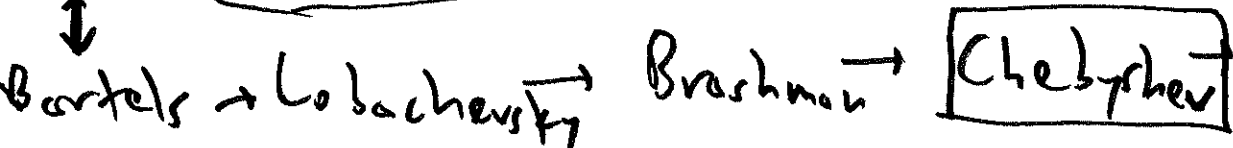
$$P(X \geq 0) = 1$$

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 Aug 17