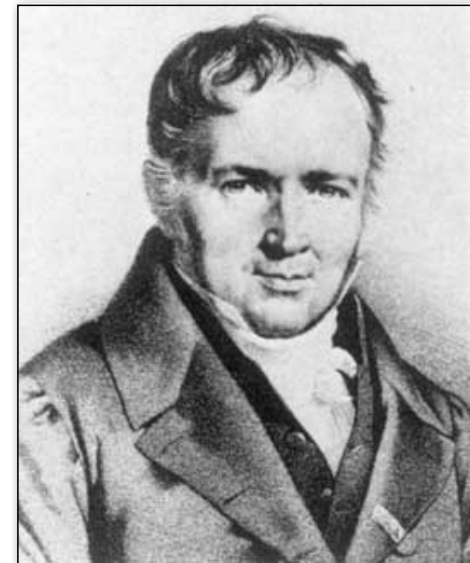
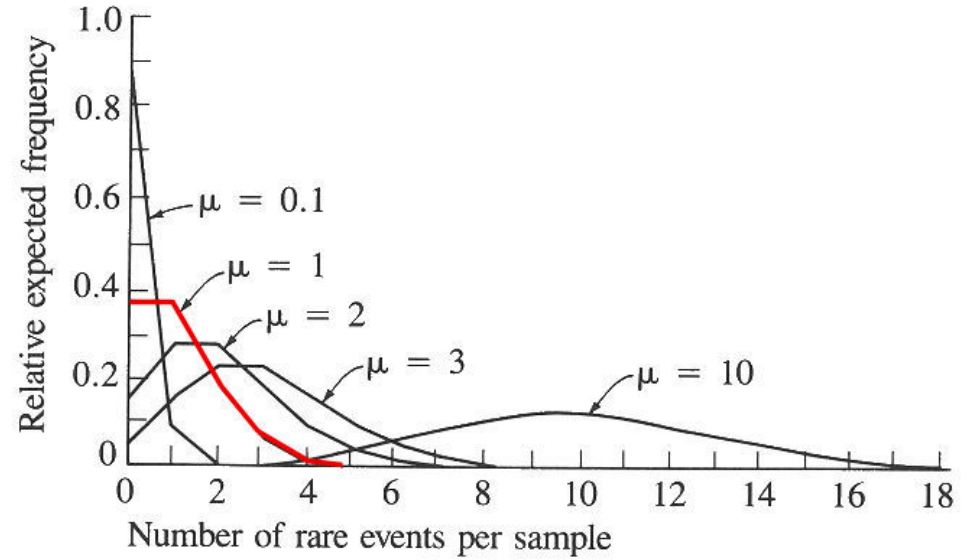
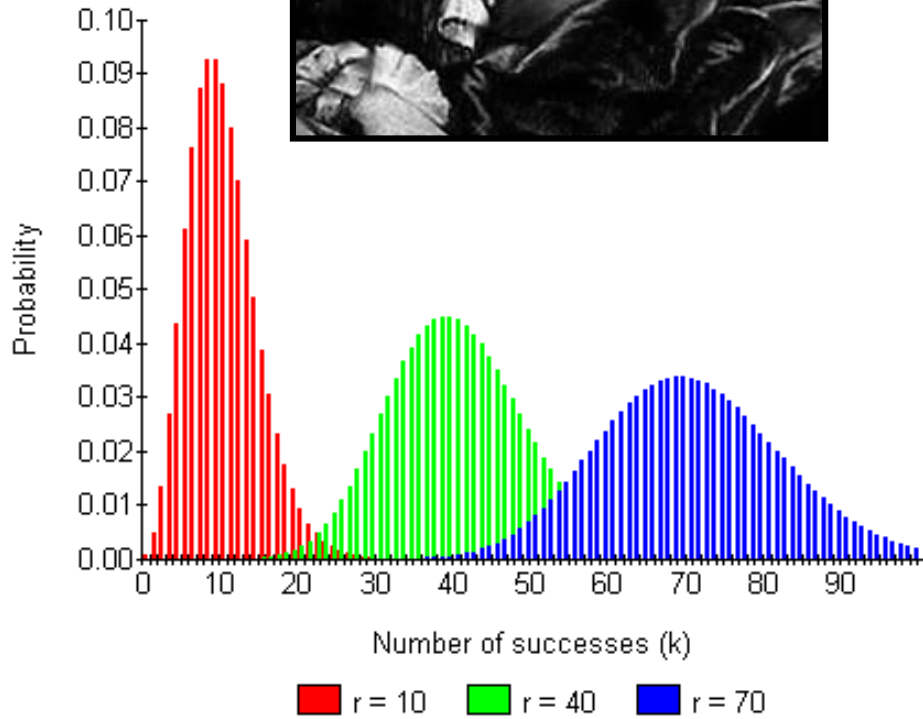


# a zoo of (discrete) random variables



Takes each possible value, say  $\{1..n\}$  with equal probability.

Say random variable “uniform on  $S$ ”

Recall envelopes problem on homework...

Randomization is key!!

Best strategy in envelopes game is randomized.

Many algorithms and protocols in CS make use of random numbers (often the best known solution)

Example: protocol for deciding when to broadcast on an Ethernet network is a randomized algorithm known as “exponential backoff”.

Consider  $n$  independent random variables  $Y_i \sim \text{Ber}(p)$

$X = \sum_i Y_i$  is the number of successes in  $n$  trials

$X$  is a *Binomial* random variable:  $X \sim \text{Bin}(n,p)$

$$P(X = i) = \binom{n}{i} p^i (1 - p)^{n-i} \quad i = 0, 1, \dots, n$$

By Binomial theorem,  $\sum_{i=0}^n P(X = i) = 1$

Examples

# of heads in  $n$  coin flips

# of 1's in a randomly generated length  $n$  bit string

# of disk drive crashes in a 1000 computer cluster

$$E[X] = pn$$

$$\text{Var}(X) = p(1-p)n$$

← (proof below, twice)

Sending a bit string over the network

$n = 4$  bits sent, each corrupted with probability  $0.1$

$X = \#$  of corrupted bits,  $X \sim \text{Bin}(4, 0.1)$

In real networks, large bit strings (length  $n \approx 10^4$ )

Corruption probability is very small:  $p \approx 10^{-6}$

Extreme  $n$  and  $p$  values arise in many cases

# bit errors in file written to disk

# of typos in a book

# of elements in particular bucket of large hash table

# of server crashes per day in giant data center

# facebook login requests sent to a particular server

## Poisson random variables

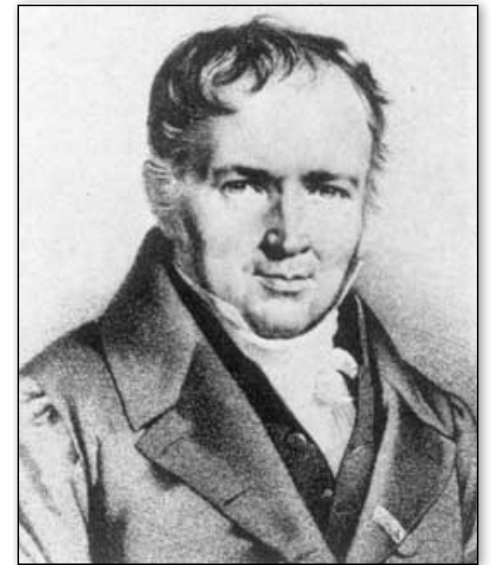
Suppose “events” happen, independently, at an *average* rate of  $\lambda$  per unit time. Let  $X$  be the *actual* number of events happening in a given time unit. Then  $X$  is a *Poisson* r.v. with *parameter*  $\lambda$  (denoted  $X \sim \text{Poi}(\lambda)$ ) and has distribution (PMF):

$$P(X = i) = e^{-\lambda} \frac{\lambda^i}{i!}$$

Examples:

- # of alpha particles emitted by a lump of radium in 1 sec.
- # of traffic accidents in Seattle in one year
- # of babies born in a day at UW Med center
- # of visitors to my web page today

See B&T Section 6.2 for more on theoretical basis for Poisson.



Siméon Poisson, 1781-1840

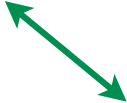
$X$  is a Poisson r.v. with parameter  $\lambda$  if it has PMF:

$$P(X = i) = e^{-\lambda} \frac{\lambda^i}{i!}$$

Is it a valid distribution? Recall Taylor series:

$$e^\lambda = \frac{\lambda^0}{0!} + \frac{\lambda^1}{1!} + \dots = \sum_{0 \leq i} \frac{\lambda^i}{i!}$$

So

$$\sum_{0 \leq i} P(X = i) = \sum_{0 \leq i} e^{-\lambda} \frac{\lambda^i}{i!} = e^{-\lambda} \sum_{0 \leq i} \frac{\lambda^i}{i!} = e^{-\lambda} e^\lambda = 1$$


## expected value of Poisson r.v.s

$$\begin{aligned} E[X] &= \sum_{0 \leq i} i \cdot e^{-\lambda} \frac{\lambda^i}{i!} && \text{ } \\ &= \sum_{1 \leq i} i \cdot e^{-\lambda} \frac{\lambda^i}{i!} && \text{ } \\ &= \lambda e^{-\lambda} \sum_{1 \leq i} \frac{\lambda^{i-1}}{(i-1)!} && \text{ } \\ &= \lambda e^{-\lambda} \sum_{0 \leq j} \frac{\lambda^j}{j!} && \text{ } \\ &= \lambda e^{-\lambda} e^{\lambda} && \text{ } \\ &= \lambda && \text{As expected, given definition} \\ & && \text{in terms of "average rate } \lambda \text{"} \end{aligned}$$

(Var[X] =  $\lambda$ , too; proof similar, see B&T example 6.20)



## binomial random variable is Poisson in the limit

---

Poisson approximates binomial when  $n$  is large,  $p$  is small, and  $\lambda = np$  is “moderate”

Formally, Binomial is Poisson in the limit as  $n \rightarrow \infty$  (equivalently,  $p \rightarrow 0$ ) while holding  $np = \lambda$

## binomial $\rightarrow$ Poisson in the limit

---

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

$$\begin{aligned} P(X = i) &= \binom{n}{i} p^i (1 - p)^{n-i} \\ &= \frac{n!}{i!(n-i)!} \left(\frac{\lambda}{n}\right)^i \left(1 - \frac{\lambda}{n}\right)^{n-i}, \text{ where } \lambda = pn \\ &= \frac{n(n-1)\cdots(n-i+1)}{n^i} \frac{\lambda^i}{i!} \frac{(1 - \lambda/n)^n}{(1 - \lambda/n)^i} \\ &= \underbrace{\frac{n(n-1)\cdots(n-i+1)}{(n-\lambda)^i}}_{\approx 1} \cdot \frac{\lambda^i}{i!} \cdot \underbrace{(1 - \lambda/n)^n}_{\approx e^{-\lambda}} \\ &\approx 1 \cdot \frac{\lambda^i}{i!} \cdot e^{-\lambda} \end{aligned}$$

I.e., Binomial  $\approx$  Poisson for large  $n$ , small  $p$ , moderate  $i$ ,  $\lambda$ .

## sending data on a network, again

---

Recall example of sending bit string over a network

Send bit string of length  $n = 10^4$

Probability of (independent) bit corruption is  $p = 10^{-6}$

$$X \sim \text{Poi}(\lambda = 10^4 \cdot 10^{-6} = 0.01)$$

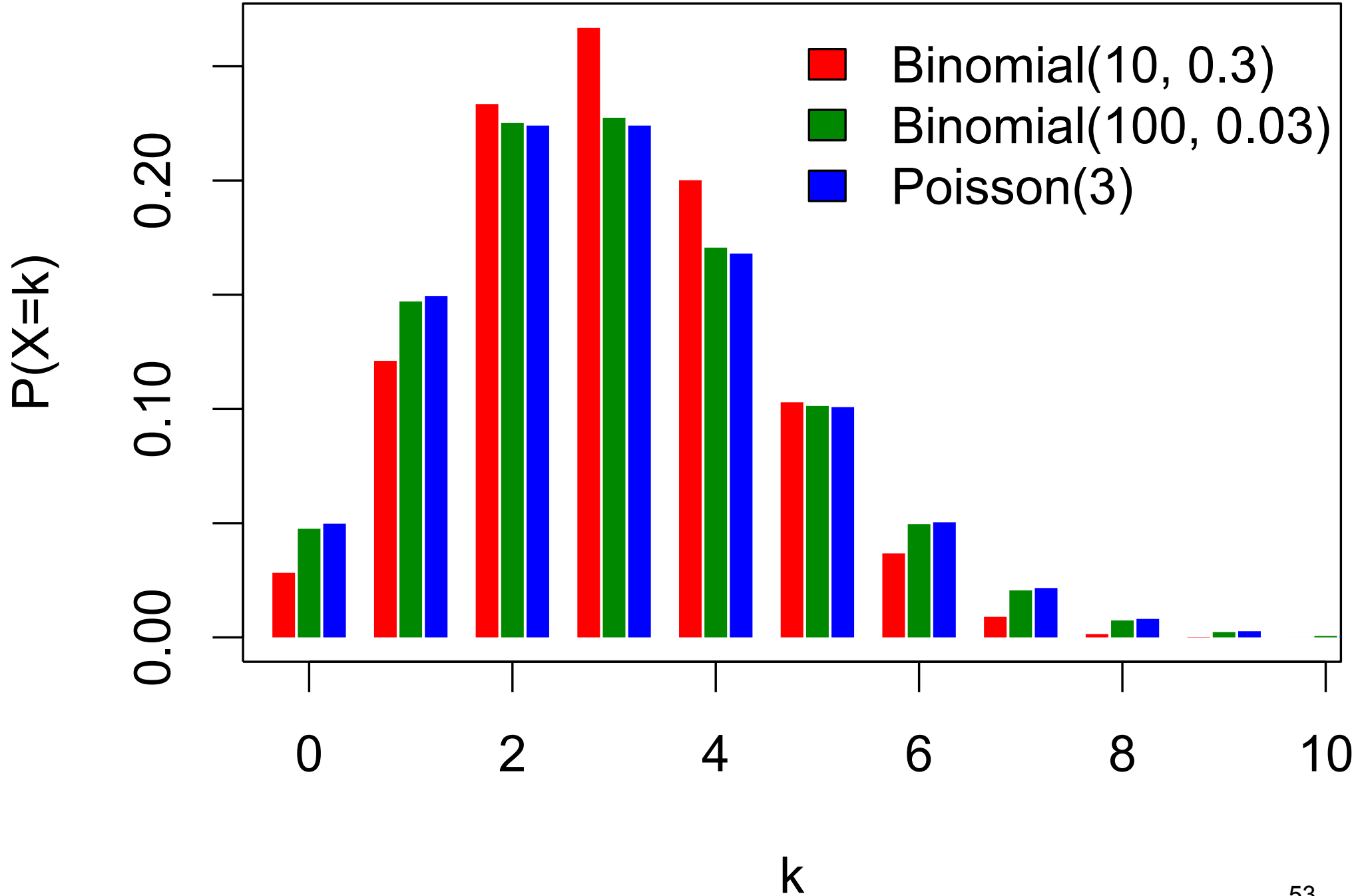
What is probability that message arrives uncorrupted?

$$P(X = 0) = e^{-\lambda} \frac{\lambda^0}{0!} = e^{-0.01} \frac{0.01^0}{0!} \approx 0.990049834$$

Using  $Y \sim \text{Bin}(10^4, 10^{-6})$ :

$$P(Y=0) \approx 0.990049829$$

## binomial vs Poisson



## expectation and variance of a poisson

---

Recall: if  $Y \sim \text{Bin}(n,p)$ , then:

$$E[Y] = np$$

$$\text{Var}[Y] = np(1-p)$$

And if  $X \sim \text{Poi}(\lambda)$  where  $\lambda = np$  ( $n \rightarrow \infty, p \rightarrow 0$ ) then

$$E[X] = \lambda = np = E[Y]$$

$$\text{Var}[X] = \lambda \approx \lambda(1-\lambda/n) = np(1-p) = \text{Var}[Y]$$

Expectation and variance of Poisson are the same ( $\lambda$ )

Expectation is the same as corresponding binomial

Variance almost the same as corresponding binomial

Note: when two different distributions share the same mean & variance, it suggests (but doesn't prove) that one may be a good approximation for the other.

In a series  $X_1, X_2, \dots$  of Bernoulli trials with success probability  $p$ , let  $Y$  be the index of the first success, i.e.,

$$X_1 = X_2 = \dots = X_{Y-1} = 0 \ \& \ X_Y = 1$$

Then  $Y$  is a *geometric* random variable with parameter  $p$ .

Examples:

Number of coin flips until first head

Number of blind guesses on SAT until I get one right

Number of darts thrown until you hit a bullseye

Number of random probes into hash table until empty slot

Number of wild guesses at a password until you hit it

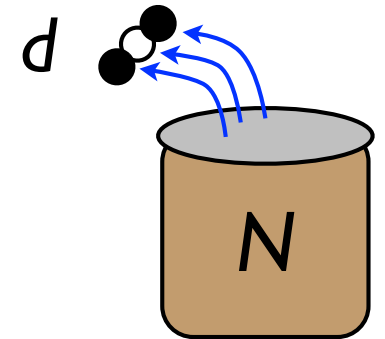
$$P(Y=k) = (1-p)^{k-1}p; \quad \text{Mean } 1/p; \quad \text{Variance } (1-p)/p^2$$

# balls in urns – the hypergeometric distribution

B&T, exercise 1.61

Draw  $d$  balls (without replacement) from an urn containing  $N$ , of which  $w$  are white, the rest black.

Let  $X$  = number of white balls drawn



$$P(X = i) = \frac{\binom{w}{i} \binom{N-w}{d-i}}{\binom{N}{d}}, \quad i = 0, 1, \dots, d$$

(note:  $\binom{n}{k} = 0$  if  $k < 0$  or  $k > n$ )

$E[X] = dp$ , where  $p = w/N$  (the fraction of white balls)

proof: Let  $X_j$  be 0/1 indicator for  $j$ -th ball is white,  $X = \sum X_j$

The  $X_j$  are *dependent*, but  $E[X] = E[\sum X_j] = \sum E[X_j] = dp$

$\text{Var}[X] = dp(1-p)(1-(d-1)/(N-1))$

$N \approx 22500$  human genes, many of unknown function

Suppose in some experiment,  $d = 1588$  of them were observed (say, they were all switched on in response to some drug)

A big question: What are they doing?

One idea: The Gene Ontology Consortium ([www.geneontology.org](http://www.geneontology.org)) has grouped genes with known functions into categories such as “muscle development” or “immune system.” Suppose 26 of your  $d$  genes fall in the “muscle development” category.

Just chance?

Or call Coach & see if he wants to dope some athletes?

Hypergeometric: GO has 116 genes in the muscle development category. If those are the white balls among 22500 in an urn, what is the probability that you would see 26 of them in 1588 draws?



**Table 2. Gene Ontology Analysis on Differentially Bound Peaks in Myoblasts versus Myotubes**

GO Categories Enriched in Genes Associated with Myotube-Increased Peaks

GOID	Term	P Value	OR <sup>a</sup>	Count <sup>b</sup>	Size <sup>c</sup>	Ont <sup>d</sup>
GO:0005856	cytoskeleton	2.05E-11	2.40	94	490	CC
GO:0043292	contractile fiber	6.98E-09	5.85	22	58	CC
GO:0030016	myofibril	1.96E-08	5.74	21	56	CC
GO:0044449	contractile fiber part	2.58E-08	5.97	20	52	CC
GO:0030017	sarcomere	4.95E-08	6.04	19	49	CC
GO:0008092	muscle tissue morphogenesis	2.50E-16	4.13	20	65	MF
GO:0007519	skeletal muscle development	2.50E-16	4.13	20	65	BP
GO:0015629	actin cytoskeleton	4.73E-06	3.08	27	111	CC
GO:0003779	actin binding	1.13E-05	2.16	27	159	MF
GO:0006936	basal body organization	3.35E-05	2.33	20	65	BP
GO:0044430	cytoskeleton part	3.35E-05	2.33	20	294	CC
GO:0031674	I band	2.27E-05	5.67	12	32	CC
GO:0003012	muscle system process	2.54E-05	4.11	16	52	BP
GO:0030029	actin filament-based process	2.89E-05	2.73	27	119	BP
GO:0007517	muscle development	5.06E-05	2.69	26	116	BP

probability of seeing this many genes from a set of this size by chance according to the hypergeometric distribution.

E.g., if you draw 1588 balls from an urn containing 490 white balls and  $\approx 22000$  black balls,  $P(94 \text{ white}) \approx 2.05 \times 10^{-11}$

A differentially bound peak was associated to the closest gene (unique Entrez ID) measured by distance to TSS within CTCF flanking domains. OR: ratio of predicted to observed number of genes within a given GO category. Count: number of genes with differentially bound peaks. Size: total number of genes for a given functional group. Ont: the Geneontology. BP = biological process, MF = molecular function, CC = cellular component.

Often care about 2 (or more) random variables *simultaneously*

measured  $X = \text{height}$  and  $Y = \text{weight}$

$X = \text{cholesterol}$  and  $Y = \text{blood pressure}$

$X_1, X_2, X_3 = \text{work loads on servers A, B, C}$

*Joint* probability mass function:

$$f_{XY}(x, y) = P(X = x \ \& \ Y = y)$$

*Joint* cumulative distribution function:

$$F_{XY}(x, y) = P(X \leq x \ \& \ Y \leq y)$$

## Two joint PMFs

W \ Z	1	2	3
1	2/24	2/24	2/24
2	2/24	2/24	2/24
3	2/24	2/24	2/24
4	2/24	2/24	2/24

X \ Y	1	2	3
1	4/24	1/24	1/24
2	0	3/24	3/24
3	0	4/24	2/24
4	4/24	0	2/24

$$P(W = Z) = 3 * 2/24 = 6/24$$

$$P(X = Y) = (4 + 3 + 2)/24 = 9/24$$

Can look at arbitrary relationships between variables this way

## Two joint PMFs

$W \backslash Z$	1	2	3	$f_W(w)$
1	2/24	2/24	2/24	6/24
2	2/24	2/24	2/24	6/24
3	2/24	2/24	2/24	6/24
4	2/24	2/24	2/24	6/24
$f_Z(z)$	8/24	8/24	8/24	

$X \backslash Y$	1	2	3	$f_X(x)$
1	4/24	1/24	1/24	6/24
2	0	3/24	3/24	6/24
3	0	4/24	2/24	6/24
4	4/24	0	2/24	6/24
$f_Y(y)$	8/24	8/24	8/24	

Marginal distribution of one r.v.:

$$f_Y(y) = \sum_x f_{XY}(x,y)$$

sum over the other:

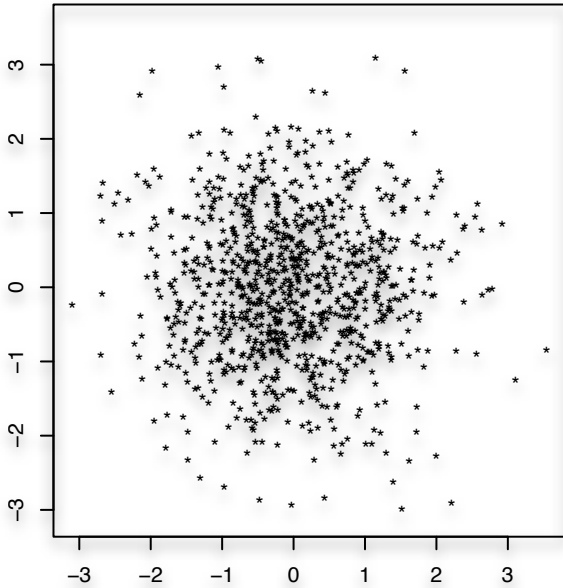
$$f_X(x) = \sum_y f_{XY}(x,y)$$

**Question:** Are W & Z independent? Are X & Y independent?

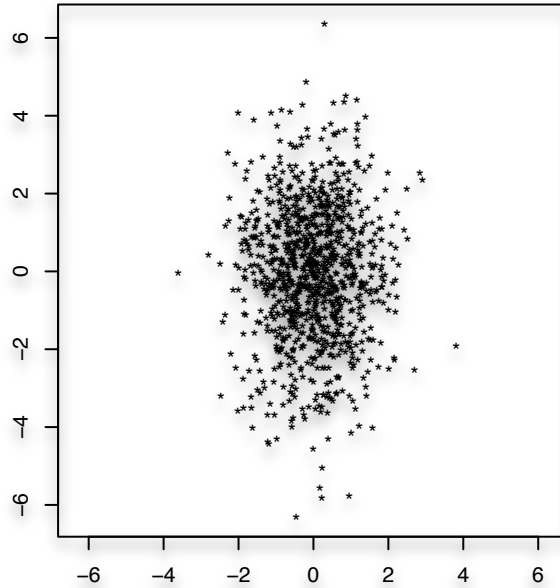
# sampling from a (continuous) joint distribution

top row: independent variables  
bottom row: dependent variables

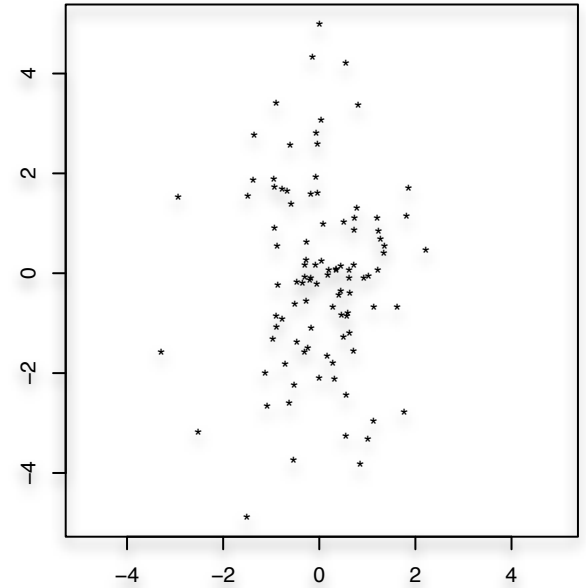
$\text{var}(x)=1, \text{var}(y)=1, \text{cov}=0, n=1000$



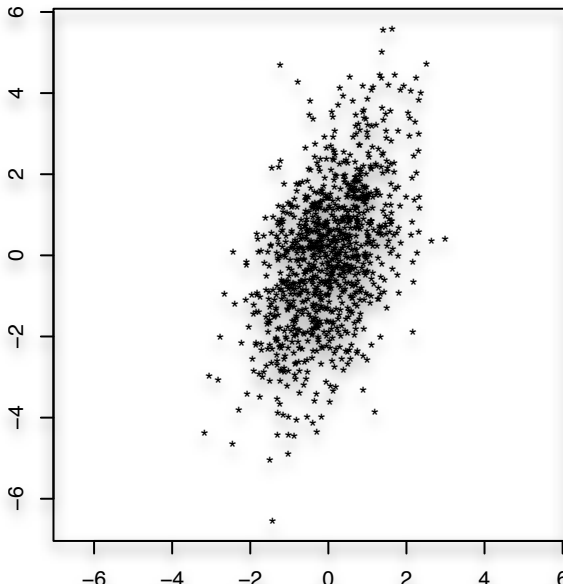
$\text{var}(x)=1, \text{var}(y)=3, \text{cov}=0, n=1000$



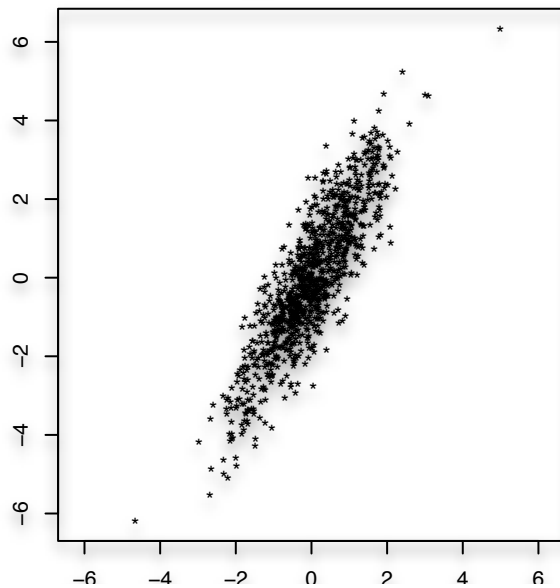
$\text{var}(x)=1, \text{var}(y)=3, \text{cov}=0, n=100$



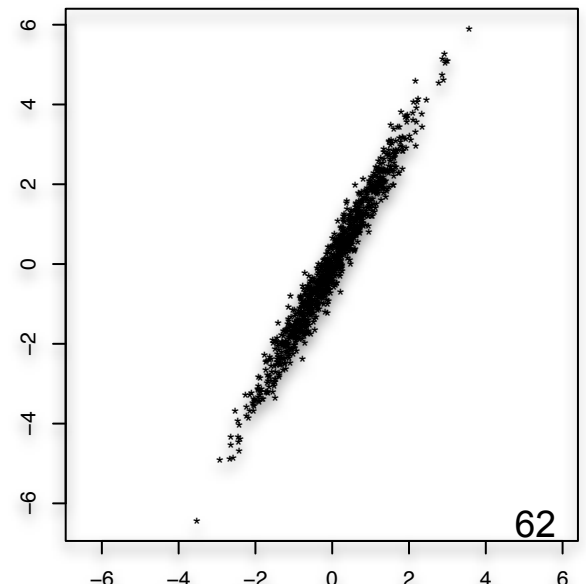
$\text{var}(x)=1, \text{var}(y)=3, \text{cov}=0.8, n=1000$



$\text{var}(x)=1, \text{var}(y)=3, \text{cov}=1.5, n=1000$



$\text{var}(x)=1, \text{var}(y)=3, \text{cov}=1.7, n=1000$



A function  $g(X, Y)$  defines a new random variable.

Its expectation is:

$$E[g(X, Y)] = \sum_x \sum_y g(x, y) f_{XY}(x, y)$$

Expectation is linear. I.e., if  $g$  is linear:

$$E[g(X, Y)] = E[aX + bY + c] = aE[X] + bE[Y] + c$$

Example:

$$g(X, Y) = 2X - Y$$

$$E[g(X, Y)] = 72/24 = 3$$

$$E[g(X, Y)] = 2 \cdot 2.5 - 2 = 3$$

X \ Y	1	2	3
1	1 • 4/24	0 • 1/24	-1 • 1/24
2	3 • 0/24	2 • 3/24	1 • 3/24
3	5 • 0/24	4 • 4/24	3 • 2/24
4	7 • 4/24	6 • 0/24	5 • 2/24

*RV*: a numeric function of the outcome of an experiment

*Probability Mass Function*  $p(x)$ : prob that  $RV = x$ ;  $\sum p(x) = 1$

*Cumulative Distribution Function*  $F(x)$ : probability that  $RV \leq x$

Concepts generalize to *joint* distributions

Expectation:

of a random variable:  $E[X] = \sum_x xp(x)$

of a function: if  $Y = g(X)$ , then  $E[Y] = \sum_x g(x)p(x)$

linearity:

$$E[aX + b] = aE[X] + b$$

$$E[X+Y] = E[X] + E[Y]; \text{ even if dependent}$$

*this interchange of “order of operations” is quite special to linear combinations. E.g.  $E[XY] \neq E[X] * E[Y]$ , in general (but see below)*

Variance:

$$\text{Var}[X] = E[(X - E[X])^2] = E[X^2] - (E[X])^2$$

Standard deviation:  $\sigma = \sqrt{\text{Var}[X]}$

$$\text{Var}[aX + b] = a^2 \text{Var}[X]$$

If  $X$  &  $Y$  are *independent*, then

$$E[X \cdot Y] = E[X] \cdot E[Y];$$

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$$

(These two equalities hold for *indp* rv's; but not in general.)



### Important Examples:

Bernoulli:  $P(X=1) = p$  and  $P(X=0) = 1-p$        $\mu = p, \sigma^2 = p(1-p)$

Binomial:  $P(X = i) = \binom{n}{i} p^i (1-p)^{n-i}$        $\mu = np, \sigma^2 = np(1-p)$

Poisson:  $P(X = i) = e^{-\lambda} \frac{\lambda^i}{i!}$        $\mu = \lambda, \sigma^2 = \lambda$

$\text{Bin}(n,p) \approx \text{Poi}(\lambda)$  where  $\lambda = np$  fixed,  $n \rightarrow \infty$  (and so  $p=\lambda/n \rightarrow 0$ )

Geometric  $P(X=k) = (1-p)^{k-1} p$        $\mu = 1/p, \sigma^2 = (1-p)/p^2$

Many others, e.g., **hypergeometric**

### Supreme Court case: Berghuis v. Smith

*If a group is underrepresented in a jury pool, how do you tell?*

Justice Breyer [Stanford Alum] opened the questioning by invoking the binomial theorem. He hypothesized a scenario involving “an urn with a thousand balls, and sixty are red, and nine hundred forty are black, and then you select them at random... twelve at a time.” According to Justice Breyer and the binomial theorem, if the red balls were black jurors then “you would expect... something like a third to a half of juries would have at least one black person” on them.

- Justice Scalia’s rejoinder: “We don’t have any urns here.”

- Should model this combinatorially
  - Ball draws not independent trials (balls not replaced)
- Exact solution:
$$P(\text{draw 12 black balls}) = \frac{\binom{940}{12}}{\binom{1000}{12}} \approx 0.4739$$
$$P(\text{draw} \geq 1 \text{ red ball}) = 1 - P(\text{draw 12 black balls}) \approx 0.5261$$
- Approximation using Binomial distribution
  - Assume  $P(\text{red ball})$  constant for every draw =  $60/1000$
  - $X = \#$  red balls drawn.  $X \sim \text{Bin}(12, 60/1000 = 0.06)$
  - $P(X \geq 1) = 1 - P(X = 0) \approx 1 - 0.4759 = 0.5240$

*In Breyer's description, should actually expect just over half of juries to have at least one black person on them*