

# Random Variables and Distributions

## Where are we going?

---

**Deer ticks:** Are they attracted by deer-gland-substance?

Suppose that 21 out of 30 deer ticks go to the deer-gland-substance-treated rod, while the other 9 go to the control rod.

- Would this be a reasonable result if the deer ticks were choosing between the rods completely at random?
- 

**Mouse survival following treatment:** Does the treatment have an effect?

Suppose that 15/30 control mice die, while 8/30 treatment mice die.

- Is the probability that a control mouse dies the same as the probability that a treatment mouse dies?

## Random variables

---

**Random variable:** A number assigned to each outcome of a random experiment.

**Example 1:** I toss a brick at my neighbor's house.

$D$  = distance the brick travels

$X$  = 1 if I break a window; 0 otherwise

$Y$  = cost of repair

$T$  = time until the police arrive

$N$  = number of people injured

**Example 2:** Treat 10 spider mites with DDT.

$X$  = number of spider mites that survive

$P$  = proportion of mites that survive.

## Further examples

---

**Example 3:** Pick a random student in the School.

$S$  = 1 if female; 0 otherwise

$H$  = his/her height

$W$  = his/her weight

$Z$  = 1 if Canadian citizen; 0 otherwise

$T$  = number of teeth he/she has

**Example 4:** Sample 20 students from the School

$H_i$  = height of student  $i$

$\bar{H}$  = mean of the 20 student heights

$S_H$  = sample SD of heights

$T_i$  = number of teeth of student  $i$

$\bar{T}$  = average number of teeth

# Random variables are ...

---

**Discrete:** Take values in a countable set  
(e.g., the positive integers).

**Example:** the number of teeth, number of gall stones, number of birds, number of cells responding to a particular antigen, number of heads in 20 tosses of a coin.

**Continuous:** Take values in an interval  
(e.g.,  $[0,1]$  or the real line).

**Example:** height, weight, mass, some measure of gene expression, blood pressure.

Random variables may also be **partly discrete and partly continuous** (for example, mass of gall stones, concentration of infecting bacteria).

## Probability function

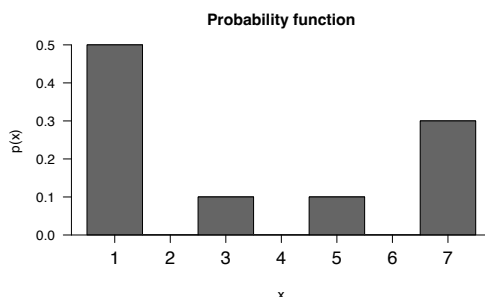
---

Consider a *discrete* random variable,  $X$ .

The **probability function** (or probability distribution, or probability mass function) of  $X$  is

$$p(x) = \Pr(X = x)$$

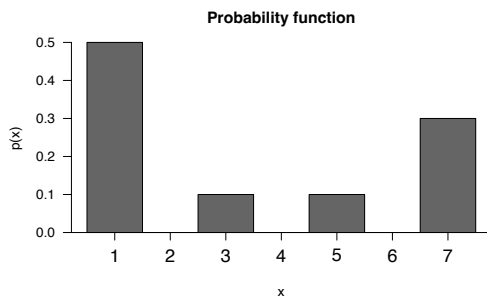
Note that  $p(x) \geq 0$  for all  $x$  and  $\sum p(x) = 1$ .



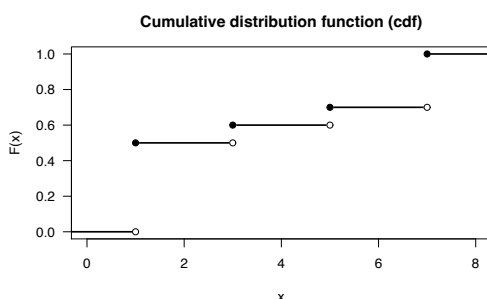
x	p(x)
1	0.5
3	0.1
5	0.1
7	0.3

# Cumulative distribution function (cdf)

The cdf of  $X$  is  $F(x) = \Pr(X \leq x)$



x	p(x)
1	0.5
3	0.1
5	0.1
7	0.3



x	F(x)
$(-\infty, 1)$	0
$[1, 3)$	0.5
$[3, 5)$	0.6
$[5, 7)$	0.7
$[7, \infty)$	1.0

# Binomial random variable

**Prototype:** The number of heads in  $n$  **independent** tosses of a coin, where  $\Pr(\text{heads}) = p$  **for each toss**.  
→  $n$  and  $p$  are called *parameters*.

Alternatively, imagine an urn containing red balls and black balls, and suppose that  $p$  is the proportion of red balls. Consider the number of red balls in  $n$  random draws *with replacement* from the urn.

**Example 1:** Sample  $n$  people at random from a large population, and consider the number of people with some property (e.g., that are graduate students or that have exactly 32 teeth).

**Example 2:** Apply a treatment to  $n$  mice and count the number of survivors (or the number that are dead).

**Example 3:** Apply a large dose of DDT to 30 groups of 10 spider mites. Count the number of groups with at least two surviving spider mites.

## Binomial distribution

---

Consider the Binomial( $n,p$ ) distribution.

That is, the number of red balls in  $n$  draws with replacement from an urn for which the proportion of red balls is  $p$ .

→ What is its probability function?

Example: Let  $X \sim \text{Binomial}(n=9,p=0.2)$ .

→ We seek  $p(x) = \Pr(X=x)$  for  $x = 0, 1, 2, \dots, 9$ .

$$p(0) = \Pr(X=0) = \Pr(\text{no red balls}) = (1-p)^n = 0.8^9 \approx 13\%.$$

$$p(9) = \Pr(X=9) = \Pr(\text{all red balls}) = p^n = 0.2^9 \approx 5 \times 10^{-7}$$

$$p(1) = \Pr(X=1) = \Pr(\text{exactly one red ball}) = \dots ?$$

## Binomial distribution

---

$$p(1) = \Pr(X=1) = \Pr(\text{exactly one red ball})$$

$$= \Pr(\text{RBBBBBBBBB or BRBBBBBBBB or ... or BBBBBBBBBR})$$

$$\begin{aligned} &= \Pr(\text{RBBBBBBBBB}) + \Pr(\text{BRBBBBBBBB}) + \Pr(\text{BBRBBBBBBB}) \\ &\quad + \Pr(\text{BBBRBBBBBB}) + \Pr(\text{BBBBRBBBBB}) \\ &\quad + \Pr(\text{BBBBBRBBBB}) + \Pr(\text{BBBBBBRBBB}) \\ &\quad + \Pr(\text{BBBBBBBBRB}) + \Pr(\text{BBBBBBBBBR}) \end{aligned}$$

$$= p(1-p)^8 + p(1-p)^8 + \dots + p(1-p)^8 = 9p(1-p)^8 \approx 30\%.$$

How about  $p(2) = \Pr(X=2)$ ?

How many outcomes have 2 red balls among the 9 balls drawn?

→ This is a problem of combinatorics. That is, counting!

## Getting at $\Pr(X=2)$

---

RRBBBBBBB RBRBBBBBBB RBBRBBBBBB RBBBRBBBBB  
RBBBBRBBB RBBBBBRBB RBBBBBBRB RBBBBBBBR  
BRRBBBBBB BRBRBBBBBB BRBBRBBBBB BRBBBRBBB  
BRBBBBRBB BRBBBBBRB BRBBBBBBR BRRRBBBBB  
BBRBRBBBB BBRBBRBBB BBRRBBRBB BBRBBBBRB  
BBRBBBBBR BBBRRBBBB BBRRRBBB BBBRBBRBB  
BBBRBBBRB BBBRBBBBR BBBBRRBBB BBBBRBRBB  
BBBBRBBRB BBBBRBBBR BBBBRRBBB BBBBRRBRB  
BBBBBRBBR BBBBBRRB BBBBBBRB BBBBBBBRR

How many are there?

$$9 \times 8 / 2 = 36.$$

## The binomial coefficient

---

The number of possible samples of size  $k$  selected from a population of size  $n$  :

$$\binom{n}{k} = \frac{n!}{k! \times (n-k)!}$$

$$\rightarrow n! = n \times (n-1) \times (n-2) \times \dots \times 3 \times 2 \times 1$$

$$\rightarrow 0! = 1$$

For a Binomial( $n,p$ ) random variable:

$$\Pr(X=k) = \binom{n}{k} p^k (1-p)^{(n-k)}$$

## Example

---

Suppose  $\Pr(\text{mouse survives treatment}) = 90\%$ , and we apply the treatment to 10 random mice.

$$\begin{aligned}\Pr(\text{ exactly 7 mice survive } ) &= \binom{10}{7} \times (0.9)^7 \times (0.1)^3 \\ &= \frac{10 \times 9 \times 8}{3 \times 2} \times (0.9)^7 \times (0.1)^3 \\ &= 120 \times (0.9)^7 \times (0.1)^3 \\ &\approx 5\%\end{aligned}$$

$$\begin{aligned}\Pr(\text{ fewer than 9 survive } ) &= 1 - p(9) - p(10) \\ &= 1 - 10 \times (0.9)^9 \times (0.1) - (0.9)^{10} \\ &\approx 26\%\end{aligned}$$

## The world is entropy driven

---

Assume we are flipping a fair coin (independently) ten times. Let  $X$  be the random variable that describes the number of heads **H** in the experiment.

$$\Pr(\text{TTTTTTTTTT}) = \Pr(\text{HTT**HH**HT**HTH**}) = (1/2)^{10}$$

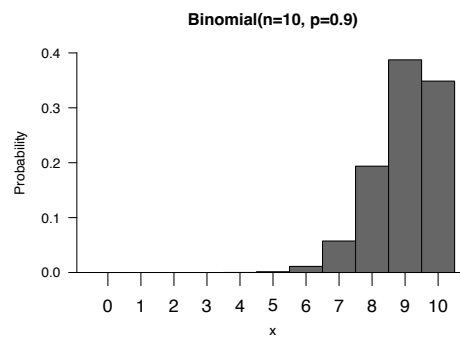
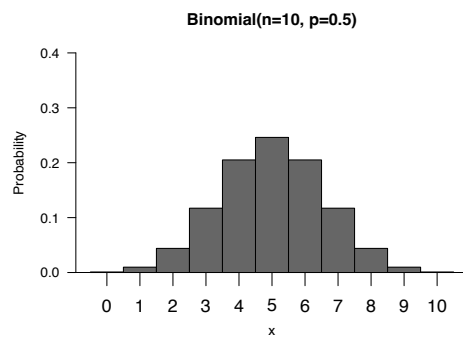
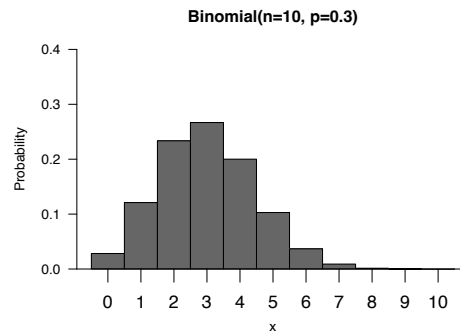
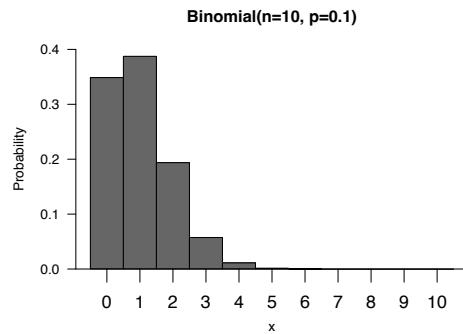
- There is only one possible outcome with zero heads.
- There are 210 possibilities for outcomes with six heads.

Thus,

- $\Pr(X = 0) = (1/2)^{10} \approx 0.1\%$ .
- $\Pr(X = 6) = 210 \times (1/2)^{10} \approx 20.5\%$ .

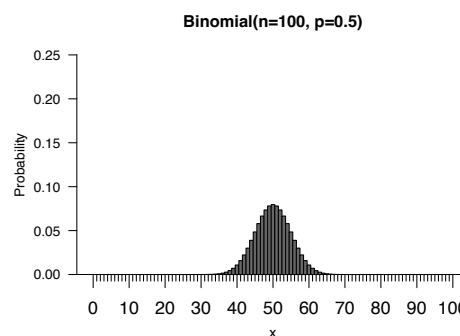
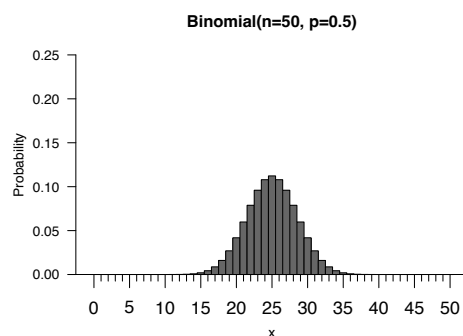
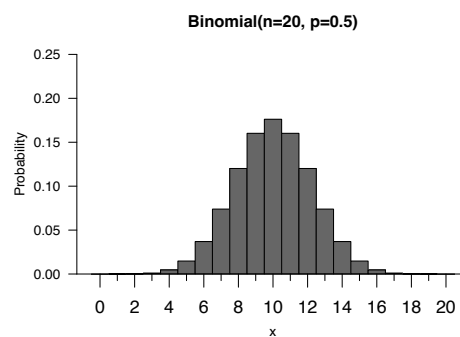
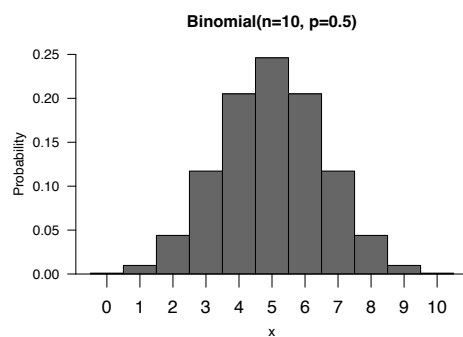
# Binomial distributions

---



# Binomial distributions

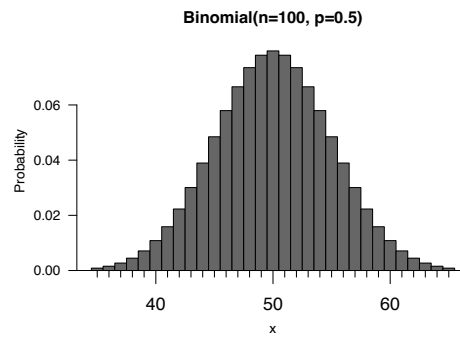
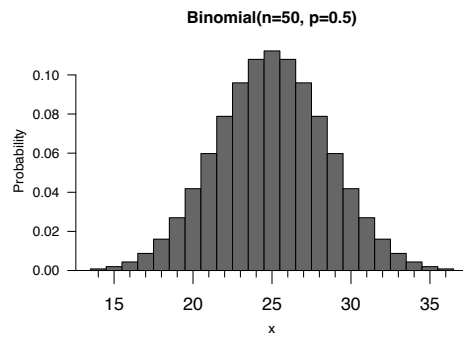
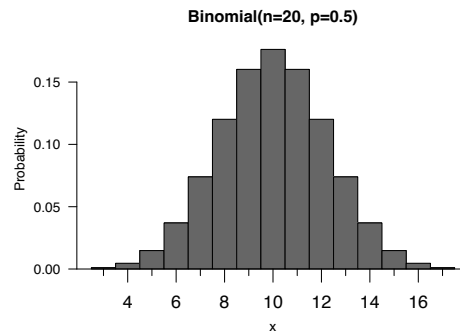
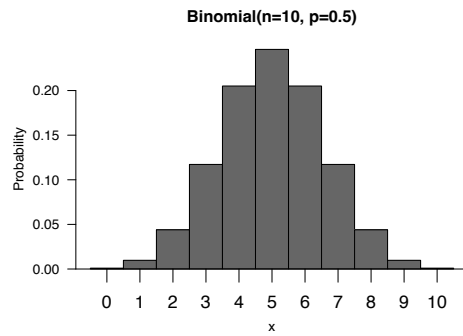
---





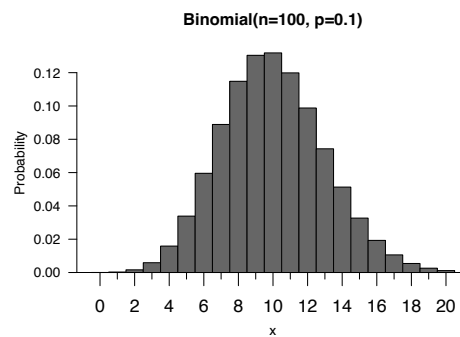
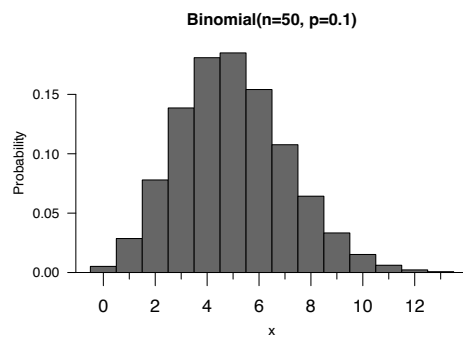
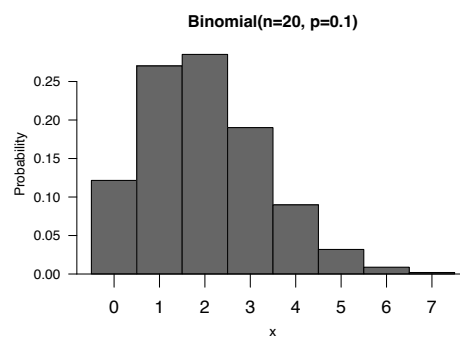
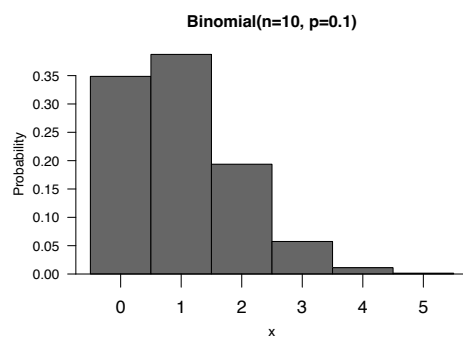
# Binomial distributions

---



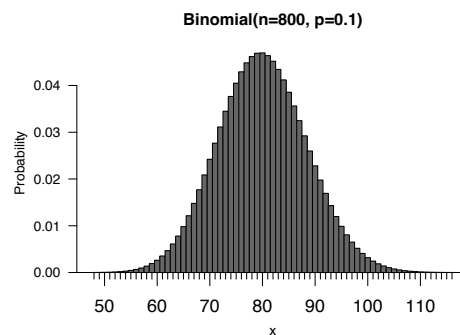
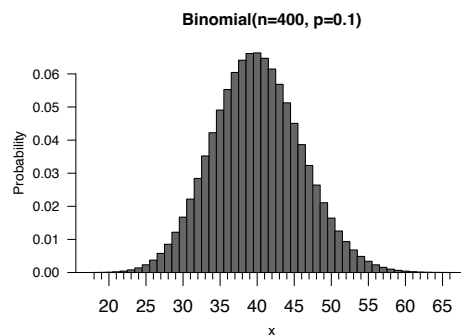
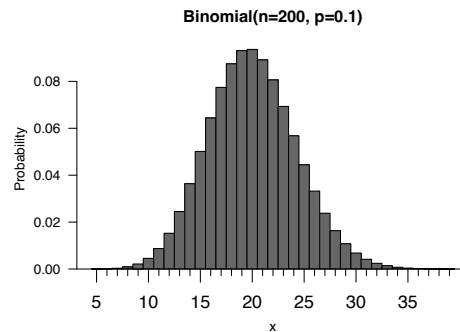
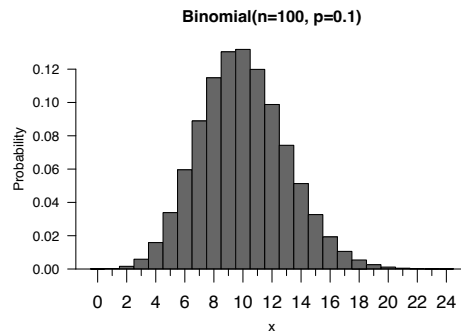
# Binomial distributions

---



# Binomial distributions

---



## Expected value and standard deviation

---

- The **expected value** (or mean) of a discrete random variable  $X$  with probability function  $p(x)$  is

$$\mu = E(X) = \sum_x x p(x)$$

- The **variance** of a discrete random variable  $X$  with probability function  $p(x)$  is

$$\sigma^2 = \text{var}(X) = \sum_x (x - \mu)^2 p(x)$$

- The **standard deviation** (SD) of  $X$  is

$$\text{SD}(X) = \sqrt{\text{var}(X)}.$$

## Mean and SD of binomial RVs

---

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

$$E(X) = n p$$

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

→ Examples:

n	p	mean	SD
10	10%	1	0.9
10	30%	3	1.4
10	50%	5	1.6
10	90%	9	0.9

## Calculations in R

---

→ Simulate binomial random variables

```
rbinom(m, size, prob)
```

→ The binomial probability function:  $\Pr(X = x)$

```
dbinom(x, size, prob)
```

→ The binomial CDF:  $\Pr(X \leq q)$

```
pbinom(q, size, prob)
```

→ The inverse CDF: the smallest  $q$  such that  $\Pr(X \leq q) \geq p$

```
qbinom(p, size, prob)
```

# Binomial random variable

---

Number of successes in  $n$  trials where:

- Trials **independent**
- $p = \text{Pr}(\text{success})$  is **constant**

The number of successes in  $n$  trials does not necessarily follow a binomial distribution!

Deviations from the binomial:

- Varying  $p$
- Clumping or repulsion (non-independence)

## Examples

---

Suppose survival differs between genders:

$\text{Pr}(\text{survive} \mid \text{male}) = 10\%$  but  $\text{Pr}(\text{survive} \mid \text{female}) = 80\%$ .

- Pick 4 male mice and 6 female mice.  
The number of survivors **is not** binomial.

- Pick 10 random mice (with  $\text{Pr}(\text{mouse is male}) = 40\%$ ).  
The number of survivors **is** binomial.

$$p = 0.4 \times 0.1 + 0.6 \times 0.8 = 0.52.$$

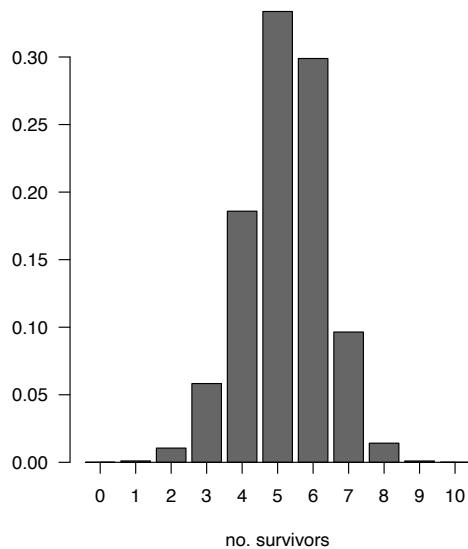
$$\text{Pr}(\text{survive}) =$$

$$\text{Pr}(\text{survive and male}) + \text{Pr}(\text{survive and female}) =$$

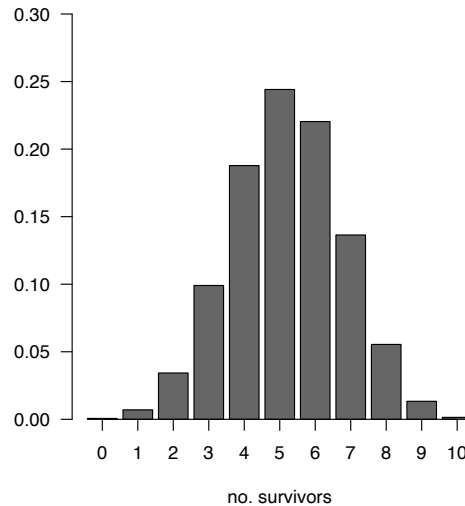
$$\text{Pr}(\text{male}) \times \text{Pr}(\text{survive} \mid \text{male}) + \text{Pr}(\text{female}) \times \text{Pr}(\text{survive} \mid \text{female})$$

# Examples

4 males; 6 females



Random mice (40% males)



# Examples

Consider Mendel's pea experiments.

Purple or white flowers, purple dominant to white:

$F_0$  genotypes are  $PP$  and  $ww$ ,  $F_1$  genotypes are  $Pw$ .

- Pick a random  $F_2$ . Self it and acquire 10 progeny.  
The number of progeny with purple flowers **is not** binomial.  
Unless we condition on the genotype of the  $F_2$  plant.
- Pick 10 random  $F_2$ 's. Self each and take a child from each.  
The number of progeny with purple flowers **is** binomial.

$$p = (1/4) \times 1 + (1/2) \times (3/4) + (1/4) \times 0 = 5/8.$$

$\Pr(\text{a progeny has a purple flower}) =$

$\Pr(\text{purple and } \{F_2 \text{ is } PP\}) + \Pr(\text{purple and } \{F_2 \text{ is } Pw\}) + \Pr(\text{purple and } \{F_2 \text{ is } ww\}) =$

$\Pr(F_2 \text{ is } PP) \times \Pr(\text{purple} | F_2 \text{ is } PP) + \Pr(F_2 \text{ is } Pw) \times \Pr(\text{purple} | F_2 \text{ is } Pw) + \Pr(F_2 \text{ is } ww) \times \Pr(\text{purple} | F_2 \text{ is } ww)$

## Multinomial distribution

---

- Imagine an urn with  $k$  types of balls.
- Let  $p_i$  denote the proportion of type  $i$ .
- Draw  $n$  balls with replacement.
- Outcome:  $(n_1, n_2, \dots, n_k)$ , with  $\sum_i n_i = n$ , where  $n_i$  is the no. balls drawn that were of type  $i$ .

$$\rightarrow P(X_1=n_1, \dots, X_k=n_k) = \frac{n!}{n_1! \times \dots \times n_k!} p_1^{n_1} \times \dots \times p_k^{n_k}$$

$$\text{if } 0 \leq n_i \leq n, \quad \sum_i n_i = n \quad \text{and} \quad 0 < p_i < 1, \quad \sum_i p_i = 1.$$

$$\text{Otherwise } P(X_1=n_1, \dots, X_k=n_k) = 0.$$

## Example

---

AA	AB	BB
35	43	22

→ Do these data correspond reasonably to the proportions 1:2:1?

Let  $(p_1, p_2, p_3) = (0.25, 0.50, 0.25)$  and  $n = 100$ .

$$P(X_1=35, X_2=43, X_3=22) = \frac{100!}{35! 43! 22!} 0.25^{35} 0.50^{43} 0.25^{22}$$
$$\approx 7.3 \times 10^{-4}$$

# Poisson distribution

---

Consider a Binomial( $n,p$ ) where

→  $n$  is really large

→  $p$  is really small

For example, suppose each well in a microtiter plate contains 50,000 T cells, and that 1/100,000 cells respond to a particular antigen.

Let  $X$  be the number of responding cells in a well.

→ In this case,  $X$  follows a **Poisson** distribution approximately.

Let  $\lambda = n p = E(X)$ .

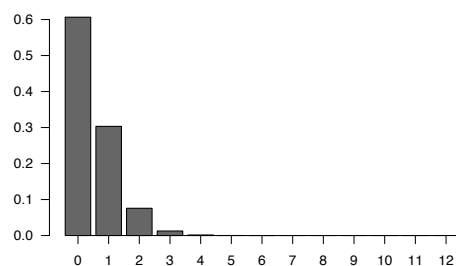
→  $p(x) = \Pr(X = x) = e^{-\lambda} \lambda^x / x!$

Note that  $SD(X) = \sqrt{\lambda}$ .

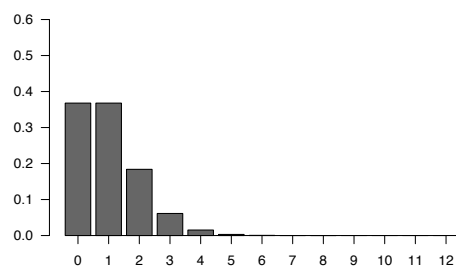
# Poisson distribution

---

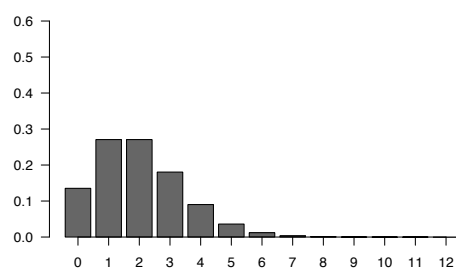
Poisson( $\lambda=1/2$ )



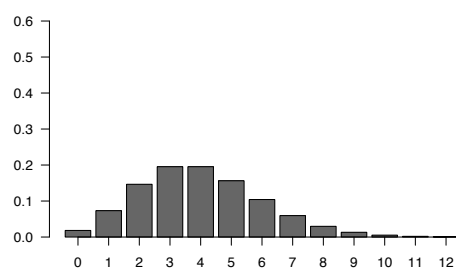
Poisson( $\lambda=1$ )



Poisson( $\lambda=2$ )



Poisson( $\lambda=4$ )



## Example

---

Suppose there are 100,000 T cells in each well of a microtiter plate. Suppose that 1/80,000 T cells respond to a particular antigen.

Let  $X$  = number of responding T cells in a well.

→  $X \sim \text{Poisson}(\lambda = 1.25)$ .

→  $E(X) = 1.25$

→  $SD(X) = \sqrt{1.25} \approx 1.12$ .

$\Pr(X = 0) = \exp(-1.25) \approx 29\%$ .

$\Pr(X > 0) = 1 - \exp(-1.25) \approx 71\%$ .

$\Pr(X = 2) = \exp(-1.25) \times (1.25)^2 / 2 \approx 22\%$ .

## Calculations in R

---

→ Simulate poisson random variables

`rpois(m, lambda)`

→ The poisson probability function:  $\Pr(X = x)$

`dpois(m, lambda)`

→ The poisson CDF:  $\Pr(X \leq q)$

`ppois(m, lambda)`

→ The inverse CDF: the smallest  $q$  such that  $\Pr(X \leq q) \geq p$

`qpois(m, lambda)`



## $Y = a + b X$

---

Suppose  $X$  is a discrete random variable with probability function  $p$ , so that  $p(x) = \Pr(X = x)$ .

→ Expected value:  $E(X) = \sum_x x p(x)$

→ Standard deviation:  $SD(X) = \sqrt{\sum_x [x - E(X)]^2 p(x)}$

Let  $Y = a + b X$  where  $a$  and  $b$  are numbers. Then  $Y$  is a random variable (like  $X$ ), and

→  $E(Y) = a + b E(X)$

→  $SD(Y) = |b| SD(X)$

In particular, if  $\mu = E(X)$ ,  $\sigma = SD(X)$ , and  $Z = (X - \mu) / \sigma$ , then

→  $E(Z) = 0$

→  $SD(Z) = 1$

## Example

---

Suppose  $X \sim \text{Binomial}(n, p)$  → number of successes

→  $E(X) = n p$

→  $SD(X) = \sqrt{n p (1 - p)}$

Let  $P = X / n$  → proportion of successes

→  $E(P) = E(X / n) = E(X) / n = p$

→  $SD(P) = SD(X / n) = SD(X) / n = \sqrt{p (1 - p) / n}$

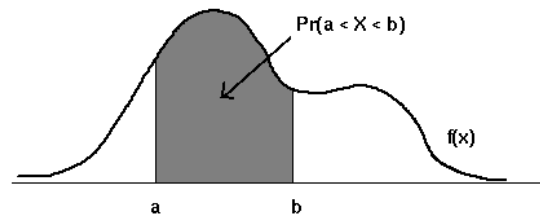
# Continuous random variables

---

Suppose  $X$  is a **continuous** random variable.

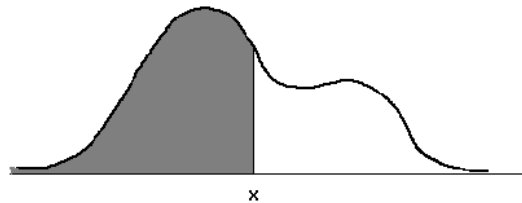
Instead of a probability function,  $X$  has a **probability density function** (pdf), sometimes called just the **density** of  $X$ .

- $f(x) \geq 0$
- $\int_{-\infty}^{\infty} f(x) dx = 1$
- Areas under curve = probabilities



Cumulative distr. function:

- $F(x) = \Pr(X \leq x)$  →



# Means and standard deviations

---

Expected value:

- Discrete RV:  $E(X) = \sum_x x p(x)$
- Continuous RV:  $E(X) = \int_{-\infty}^{\infty} x f(x) dx$

Standard deviation:

- Discrete RV:  $SD(X) = \sqrt{\sum_x [x - E(X)]^2 p(x)}$
- Continuous RV:  $SD(X) = \sqrt{\int_{-\infty}^{\infty} [x - E(X)]^2 f(x) dx}$

# Uniform distribution

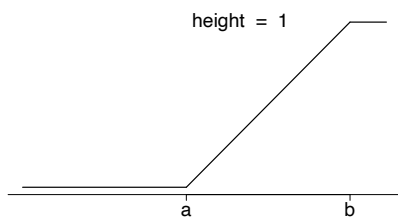
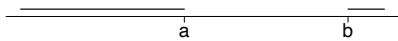
---

$X \sim \text{Uniform}(a, b)$

→ Draw a number at random from the interval  $(a, b)$ .

height =  $\frac{1}{b-a}$  \_\_\_\_\_

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a < x < b \\ 0 & \text{otherwise} \end{cases}$$



→  $E(X) = (a + b) / 2$

→  $SD(X) = (b - a) / \sqrt{12}$   
 $\approx 0.29 \times (b - a)$

# Normal distribution

---

By far the most important distribution:

The **normal distribution** (also called the Gaussian distribution).

If  $X \sim N(\mu, \sigma)$ , then the pdf of  $X$  is

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

Note:  $E(X) = \mu$  and  $SD(X) = \sigma$ .

Of great importance:

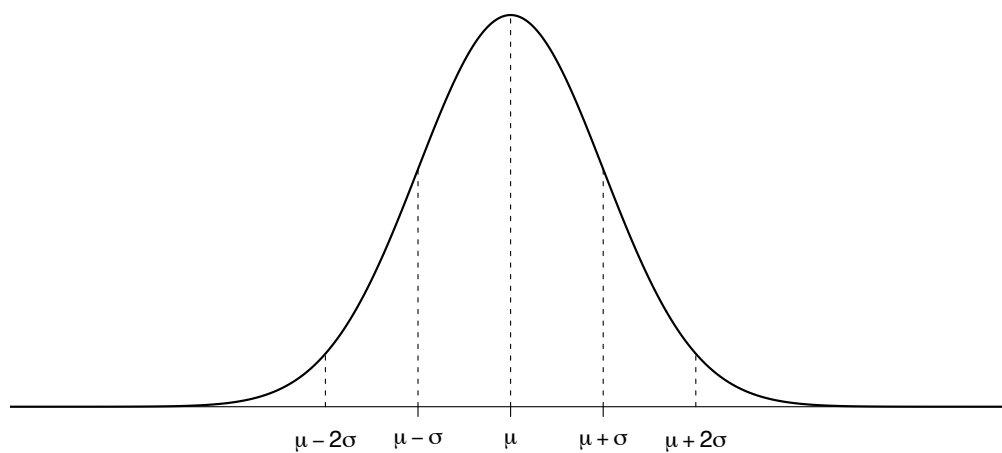
→ If  $X \sim N(\mu, \sigma)$  and  $Z = (X - \mu) / \sigma$ , then  $Z \sim N(0, 1)$ .

This is the **standard normal distribution**.

# Normal distribution

---

The normal curve

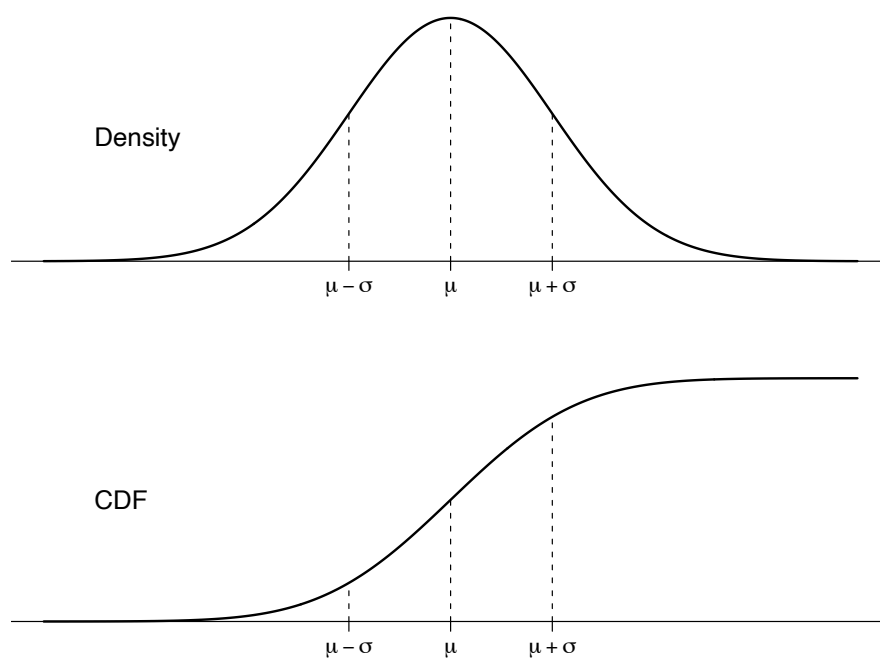


→ Remember:

$$\Pr(\mu - \sigma \leq X \leq \mu + \sigma) \approx 68\% \quad \text{and} \quad \Pr(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \approx 95\%.$$

# The normal CDF

---

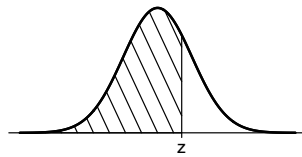


# Calculations with the normal curve in R

---

- Convert to a statement involving the cdf.
- Use the function `pnorm()`.

→ Draw a picture!

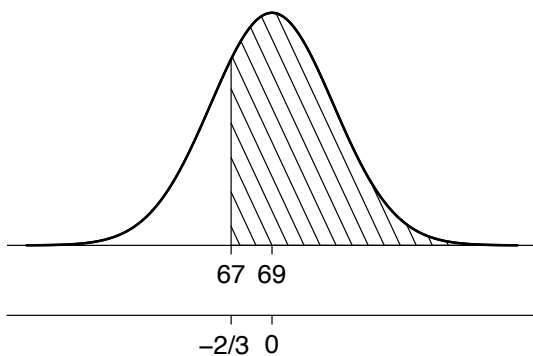


## Examples

---

Suppose the heights of adult males in the U.S. are approximately normal distributed, with mean = 69 in and SD = 3 in.

→ What proportion of men are taller than 5'7"?



$$X \sim N(\mu=69, \sigma=3)$$

$$Z = (X - 69)/3 \sim N(0,1)$$

$$\Pr(X \geq 67) =$$

$$\Pr(Z \geq (67 - 69)/3) =$$

$$\Pr(Z \geq -2/3)$$

# R

---



Use either of the following three:

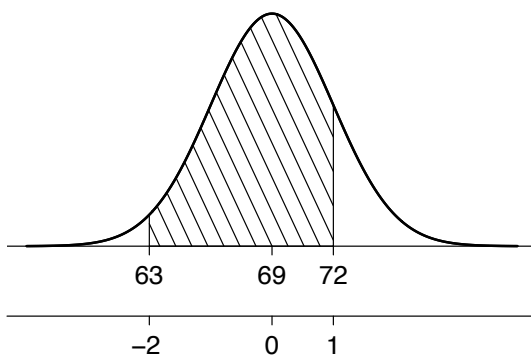
- `pnorm(67, 69, 3, lower=FALSE)`
- `1 - pnorm(67, 69, 3)`
- `pnorm(2/3)`

The answer: 75%.

## Another calculation

---

- What proportion of men are between 5'3" and 6'?

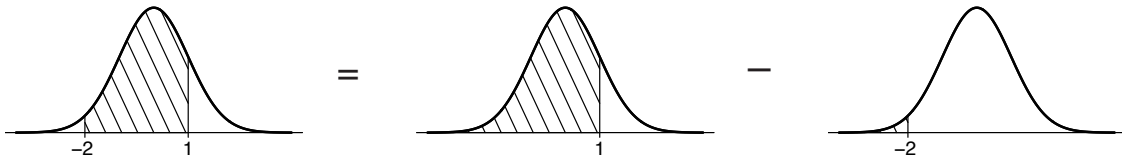


$$\Pr(63 \leq X \leq 72) =$$

$$\Pr(-2 \leq Z \leq 1)$$

# R

---



Use either of the following:

→ `pnorm(72, 69, 3) - pnorm(63, 69, 3)`

→ `pnorm(1) - pnorm(-2)`

The answer: 82%.