

## Lecture 25: Using the Normal Approximation: Ross 5.4 (Mendel's experiments)

### 25.1: Using the Normal probability table

The table in Ross, P.222, is of the usual form for the probabilities for a  $N(0,1)$  standard Normal distribution. It gives  $P(Z \leq x)$  for values of  $x$  from 0 up. This probability is denoted  $\Phi(x)$ .

For negative  $x$ ,  $P(Z \leq x) = P(Z \geq -x) = 1 - P(Z \leq -x)$ .

(Note that since  $Z$  is a continuous random variable,  $P(Z < x) = P(Z \leq x)$ ).

For general  $a, b$ :  $P(a < Z \leq b) = \Phi(b) - \Phi(a)$ .

### 25.2: Mendel's experiments

Mendel did many experiments of the form of the one with the red/white flowers. He crossed red-flowered plants with white-flowered plants, so he knew the red-flowered offspring were of RW type. These are known as the  $F_1$  or *hybrids*. He then crossed these with each other, and expected to get red and white flowers in the ratio 3:1.

Here are four examples:

- a) 253  $F_1$  producing 7324 seeds: 5474 round, 1850 wrinkled: ratio 2.96:1
- b) 258  $F_1$  producing 8023 seeds: 6022 yellow, 2001 green: ratio 3.01:1.
- c) 929  $F_2$ ; 705 red flowers, 224 white flowers: ratio 3.15:1.
- d) 580  $F_2$ : 428 green pods, 152 yellow pods: ratio 2.82:1

### 25.3 Are Mendel's results too good?

There has been much debate as to whether Mendel's results are "too good" – too close to the 3:1 ratio.

Note the larger samples for characteristics that can be observed at the seed stage. These give the ratios closest to 3:1. This is as expected:  $\text{var}(X) = np(1-p)$  but  $\text{var}(X/n) = \text{var}(X)/n^2 = p(1-p)/n$  which decreases as  $n$  increases. Are we too close? Recall  $Z = (X - np)/\sqrt{np(1-p)}$  is approx  $N(0,1)$ . Here  $p = 3/4$ :

- a)  $Z_a = (5474 - 7324 \times 0.75)/\sqrt{7324 \times 3/16} = -0.5127$ ,  $P(-0.5127 < Z \leq 0.5127) = 2\Phi(0.5127) - 1 = 0.39$ .
- b)  $Z_b = (6022 - 8023 \times 0.75)/\sqrt{8023 \times 3/16} = 0.1225$ ,  $P(-0.1225 < Z \leq 0.1225) = 2\Phi(0.1225) - 1 = 0.097$ .
- c)  $Z_c = (705 - 929 \times 0.75)/\sqrt{929 \times 3/16} = 0.6251$ ,  $P(-0.6251 < Z \leq 0.6251) = 2\Phi(0.6251) - 1 = 0.468$ .
- d)  $Z_d = (428 - 580 \times 0.75)/\sqrt{580 \times 3/16} = -0.6712$ ,  $P(-0.6712 < Z \leq 0.6712) = 2\Phi(0.6712) - 1 = 0.498$ .

So far, with these experiments, there seems no reason to think Mendel's results are "too good".

### 25.4 Combining the experiments

The fact that these involve different characteristics does not stop us combining them. They are all independent Bernoulli trials with  $p = 0.75$ .

We have  $7324 + 8023 + 929 + 580 = 16856$  trials with  $5474 + 6022 + 705 + 428 = 12629$  "successes".  $Z = (12629 - 16856 \times 0.75)/\sqrt{16856 \times 3/16} = -0.2312$ .  $P(-0.2312 < Z \leq 0.2312) = 2\Phi(0.2312) - 1 = 0.183$ .

Alternatively, we can combine the  $Z$ -values: we could do this even if they came from Bernoulli trials with different  $p$ .

Here:  $Z_a + Z_b + Z_c + Z_d = -0.5127 + 0.1225 + 0.6251 - 0.6712 = -0.4363$ .

This would be a Normal with mean 0 but variance 4 (why?). So we must standardize it:

$Z^* = -0.4363/2 = -0.2182$ ,  $P(-0.2182 < Z \leq 0.2182) = 2\Phi(0.2182) - 1 = 0.173$ .

So again, either way, here there is no evidence of the results being "too good".

However, when a large number of Mendel's other results are also grouped together, overall, they do look a bit "too good".

## Lecture 26: More examples from Mendel's experiments

### 26.1 Approximating discrete Binomial with continuous Normal

(a) When approximating  $P(X = k)$  for a binomial  $X$  by a Normal  $Y$ , strictly we should consider  $P(k - \frac{1}{2} < Y \leq k + 1/2)$  (see homework example).

(b) However, for large  $n$  it makes almost no difference. Recall when  $X$  is increased by 1,  $Z$  is increased by  $\delta = 1/\sqrt{np(1-p)}$ .

(c) Example: suppose  $X$  is  $Bin(30, 2/3)$ .  $E(X) = 20$ ,  $\text{var}(X) = 30 \times (1/3) \times (2/3) = 20/3$ .

Compute the probability  $14 \leq X \leq 18$

(i) Exactly, using the Binomial probabilities:  $\sum_{k=14}^{18} P(X = k)$ . Answer 0.2689.

(ii) Using the Normal approx, with the range 14 to 18 for  $X$ :

$$Z = (14 - 20)/\sqrt{20/3} = -2.32 \text{ to } Z = (18 - 20)/\sqrt{20/3} = -0.77. \text{ Answer: } 0.2105.$$

(iii) Using the Normal approx, with the range 13.5 to 18.5 for  $X$ :

$$Z = (13.5 - 20)/\sqrt{20/3} = -2.517 \text{ to } Z = (18.5 - 20)/\sqrt{20/3} = -0.5809. \text{ Answer: } 0.2747.$$

### 26.2 Mendel's experiment: continued

Now Mendel wanted to show not just the 3:1 red:white ratio, but also the 1:2:1 for  $RR : RW : WW$ . So he needed to find which of his red-flowered  $F_2$  plants were  $RR$  and which were  $RW$ . To do this he *selfed* his red-flowered  $F_2$  pea plants: that is, the parents were  $RR$  giving  $RR \times RR$  or  $RW$  giving  $RW \times RW$ .

In order to tell whether the parent was  $RW$ , Mendel grew up 10 offspring, and if all were red he said the plant *bred true*. Note, under Mendel's hypothesis  $P(RR \mid \text{red}) = 1/3$ .

Mendel reported his result: from 600  $F_2$  he found 201 *bred true*. Assuming  $1/3$  should *breed true*, is this result too close to  $1/3$ ? Note if  $p = 1/3$ ,  $E(X) = 200$ ,  $\text{var}(X) = 600 \times 1/3 \times 2/3 = 400/3$ .

(i) Without the correction (considering  $X = 199, 200, 201$ ) show the probability of being this close is about 6.5%. ( $Z = \pm 0.08660$ ).

(ii) With the correction ( $189.5 < X < 201.5$ ) show the probability of being this close is a bit over 10% ( $Z = \pm 0.12990$ ).

(Here the continuity correction makes enough difference that it might affect our belief about whether Mendel's results are "too good").

### 26.3 Mendel's mistake:

Recall that each offspring of an  $RW \times RW$  mating is white with probability  $1/4$ .

(i) For each  $RW \times RW$  mating, what is the probability Mendel mis-called it as  $RR \times RR$ ?

$$\text{Answer: } (3/4)^{10} = 0.0563.$$

(ii) If the frequency of  $RR$  parents is  $1/3$  and  $RW$  is  $2/3$ , what is the overall probability that all 10 offspring plants are red? Answer:  $(1/3) + (2/3) \times 0.0563 = 0.371$ .

### 26.4 Probability of being close to 0.371

So now the  $p$  of Mendel's Binomial should have been  $p = 0.371$ .  $E(X) = 222.6$ ,  $\text{var}(X) = 140.01$ ,  $\text{st.dev} = 11.83$ . Now we need the probability that Mendel's reported count of 201 would be *this far off*.

(i) With no correction:  $X \leq 201$ ,  $Z < -1.825$  or  $Z > 1.825$ . Answer: about 6.8%.

(ii) With correction:  $X \leq 201.5$ ,  $Z < -1.783$  or  $Z > 1.783$ . Answer: about 7.4%.

(iii) Or maybe we should ask, this far off in direction of his assumed  $1/3$ , Answers: 3.4% and 3.7%.

Either Mendel was, for once, quite *unlucky* or else his result is too close to what he may have expected, and too far from what he should have found.

## Lecture 27: The Cumulative Distribution Function: Ross 4.9, 5.2

**27.1 (i) Definition:** (Ross 4.1) For any random variable  $X$ , the *cumulative distribution function* is defined as  $F_X(x) = P(X \leq x)$  for  $-\infty < x < \infty$ .

(ii) For a discrete random variable with pmf  $p_X(x)$ ,  $F_X(b) = \sum_{x \leq b} p_X(x)$ .

(iii) For a continuous random variable with pdf  $f_X(x)$ ,  $F_X(b) = \int_{-\infty}^b f_X(x) dx$ .

(iv) For all random variables,  $P(a < X \leq b) = F(b) - F(a)$

because  $\{X \leq b\} = \{X \leq a\} \cup \{a < X \leq b\}$  and  $\{X \leq a\} \cap \{a < X \leq b\} = \Phi$  (empty set).

### 27.2 Properties: (Ross 4.9)

(i)  $F_X$  is a non-decreasing function: if  $a < b$ , then  $F_X(a) \leq F_X(b)$ , because  $\{X \leq a\} \subset \{X \leq b\}$ .

(ii)  $\lim_{b \rightarrow \infty} F_X(b) = 1$ , because for any increasing sequence  $b_n \rightarrow \infty$ ,  $n = 1, 2, 3, \dots$ ,

$\Omega = \{X < \infty\} = \cup \{X \leq b_n\}$ , so  $1 = P(\Omega) = \lim_{n \rightarrow \infty} P(X \leq b_n) = \lim_{n \rightarrow \infty} F_X(b_n)$ .

(iii)  $\lim_{b \rightarrow -\infty} F_X(b) = 0$ , because for any decreasing sequence  $b_n \rightarrow -\infty$ ,  $n = 1, 2, 3, \dots$ ,

$\Phi = \{X = -\infty\} = \cap \{X \leq b_n\}$ , so  $0 = P(\Phi) = \lim_{n \rightarrow \infty} P(X \leq b_n) = \lim_{n \rightarrow \infty} F_X(b_n)$ .

(iv)  $F_X$  is right-continuous. That is, for any  $b$  and any decreasing sequence  $b_n$ ,  $n = 1, 2, 3, \dots$ , with  $b_n \rightarrow b$  as  $n \rightarrow \infty$ ,  $\lim_{n \rightarrow \infty} F_X(b_n) = F_X(b)$ , because  $\{X \leq b\} = \cap \{X \leq b_n\}$ .

Note  $P(X \leq b) = P(X < b) + P(X = b)$ , and  $P(X < b) = \lim_{x \rightarrow b^-} F(x)$ .

If  $X$  is discrete, with  $P(X = b) > 0$ ,  $F_X$  will be discontinuous at  $x = b$ .

### 27.3 Case of continuous random variables: (Ross 5.2)

For discrete random variables,  $F_X(x)$  is just a set of flat (constant) pieces, with jumps in amount  $P(X = x_i)$  at each possible value  $x_i$  of  $X$ . This is not very useful.

For continuous random variables, the cdf is very useful!

$$F_X(x) = P(X \leq x) = \int_{-\infty}^x f_X(w) dw \quad \text{so} \quad \frac{dF_X(x)}{dx} = f_X(x).$$

That is, we get the pdf by differentiating the cdf: the cdf is often easier to consider.

Example: scaling an exponential random variable.

Suppose  $f_X(x) = \lambda e^{-\lambda x}$  on  $x \geq 0$ , and let  $Y = aX$  ( $a > 0$ ). What is the pdf of  $Y$ ?

$$\text{First, } F_X(x) = \int_0^x \lambda e^{-\lambda w} dw = [-e^{-\lambda w}]_0^x = 1 - e^{-\lambda x} \text{ on } x \geq 0.$$

$$\text{Now, } F_Y(y) = P(Y \leq y) = P(aX \leq y) = P(X \leq y/a) = F_X(y/a) = (1 - e^{-\lambda y/a}),$$

$$\text{so } f_Y(y) = F'_Y(y) = \frac{d}{dy}(1 - e^{-\lambda y/a}) = (\lambda/a)e^{-(\lambda/a)y} \text{ on } y \geq 0.$$

That is  $Y$  is an exponential random variable with parameter  $\lambda/a$ .

### 27.4 Using the cdf to consider functions of random variables

Using the cdf is often the easiest way to consider functions of a random variable.

Example: Suppose  $X$  is Uniform  $U(0,1)$ . What is the pdf of  $Y = X^3$ ?

$$f_X(x) = 1, \quad 0 \leq x \leq 1; \quad F_X(x) = x, \quad 0 \leq x \leq 1$$

$$F_Y(y) = P(Y \leq y) = P(X^3 \leq y) = P(X \leq y^{1/3}) = F_X(y^{1/3}) = y^{1/3}, \quad 0 \leq y \leq 1$$

$$f_Y(y) = \frac{d}{dy} F_Y(y) = (1/3)y^{-2/3} \quad 0 \leq y \leq 1$$

$$\text{Note: } E(X^3) = \int_0^1 x^3 dx = 1/4. \quad E(Y) = \int_0^1 y(1/3)y^{-2/3} dy = [(1/3)y^{4/3}/(4/3)]_0^1 = 1/4$$