The solutions which we reproduce below are far more elaborate than what was expected on the final. After all, not many students have laser printers and Mathematica running on a laptop during the exam! Our purpose is to give as full an explanation as possible, so *you* understand the solution.

*Problem* 1. (20 points) In a company of 200 employees, there are 32 employees making at least $100,000 a year. There are 47 employees in the company that have a graduate degree. There are 143 employees that do not have a graduate degree and earn less than $100,000 per year.

  (a) Find the probability that a randomly chosen employee makes less than $100,000 and has a graduate degree.

  (b) Find the probability that a randomly chosen employee makes at least $100,000 or has a graduate degree.

  (c) If an employee is selected at random, let $A$ be the event that the employee makes less than $100,000, and let $B$ be the event that the employee does not have a graduate degree. Find the value of $P(B|A)$.

  (d) Are the events $A$ and $B$ in part (c) independent? Justify your answer.

  (e) Two employees are selected randomly to attend a lunch with the CEO of the company. What is the probability they both have a graduate degree?

*Solution.* The easiest way to analyze this problem is to write down a table summarizing the information which we know. We'll use two columns, one denoted $G$ (for those with graduate degrees) and one denoted $\overline{G}$ (for those *without* graduate degrees; and two rows, one labeled $R$ (for 'Rich', er, those making $100,000 or more) and one labeled $\overline{R}$ (for those making less than $100,000). We'll mark the row sums to the right of the respective rows, and the column sums just below the respective columns. In the lower right column we enter 200, which should be both the row sum and the column sum for that column.

Here's the result:

|  | $G$ | $\overline{G}$ | Total |
|---|---|---|---|
| $R$ |  |  | 32 |
| $\overline{R}$ |  | 143 |  |
| Total | 47 |  | 200 |

The first thing we notice is that in order for the last row to sum to 200, the column sum for $\overline{G}$ must be 153; and in order for the last column to sum to 200, the entry opposite $\overline{R}$ must be 168. So here's what we know:

|  | $G$ | $\overline{G}$ | Total |
|---|---|---|---|
| $R$ |  |  | 32 |
| $\overline{R}$ |  | 143 | 168 |
| Total | 47 | 153 | 200 |

But this immediately allows us to fill in the $\overline{R} \cap G$ entry ($25 = 168 - 143$) and the $R \cap \overline{G}$ entry ($10 = 153 - 143$):

|  | $G$ | $\overline{G}$ | Total |
|---|---|---|---|
| $R$ |  | 10 | 32 |
| $\overline{R}$ | 25 | 143 | 168 |
| Total | 47 | 153 | 200 |

And finally we fill in the $R \cap G$ entry either from the row ($22 = 32 - 10$) *or* from the column ($22 = 47 - 25$):

|  | $G$ | $\overline{G}$ | Total |
|---|---|---|---|
| $R$ | 22 | 10 | 32 |
| $\overline{R}$ | 25 | 143 | 168 |
| Total | 47 | 153 | 200 |

Now we're ready to answer the questions.

  (a) "Makes less than $100,000" is $\overline{R}$; "has a graduate degree" is $G$; there are 25 people in $\overline{R} \cap G$, as we can see from the final filled-in table, so the probability of randomly choosing one of these 25 from the 200 total employees is $25/200$, or 0.125.

  (b) "Makes at least $100,000" is $R$; "has a graduate degree" is $G$. The union of column $G$ and row $R$ contains $22 + 10 + 25 = 57$ individuals. Be careful that you don't just add $47 + 32$, the number in column $G$ and row $R$ because this will double-count the individuals who are in $R \cap G$. So the probability is $57/200 = 0.285$.

  (c) By definition, $P(B|A) = \dfrac{P(B \cap A)}{P(A)}$. Now event $A$ is line $\overline{R}$ of the table, event $B$ is column $\overline{G}$. Thus $B \cap A$ contains 143 individuals, thus $P(B \cap A) = 143/200$; and $P(A) = 168/200$; thus

$$P(B|A) = \frac{143/200}{168/200} = \frac{143}{168} = 0.8512.$$

(d) In order for $B$ and $A$ to be independent we must have $P(B|A) = P(B)$. But $P(B|A) = 0.8512$ and $P(B) = 153/200 = 0.765$, and these are *not* equal; thus $A$ and $B$ are *not* independent.

(e) The probability is

$$\frac{47}{200} \times \frac{46}{199} = 0.0543216.$$

It will be a dull lunch. Note that this is selection *without replacement*; having chosen one person with a graduate degree for lunch, the CEO has reduced both the number of employees from which to pick, and the number of employees with graduate degrees.

It's worth remarking in problems like this (a $2 \times 2$ table) that we can think of the problem as consisting of nine unknowns, call them

$$x_1, x_2, x_3, y_1, y_2, y_3, z_1, z_2, z_3,$$

which are the entries of a table:

|  | $G$ | $\overline{G}$ | Total |
|---|---|---|---|
| $R$ | $x_1$ | $y_1$ | $z_1$ |
| $\overline{R}$ | $x_2$ | $y_2$ | $z_2$ |
| Total | $x_3$ | $y_3$ | $z_3$ |

connected by six equations representing the fact that some of the entries are row or column sums:

$$x_1 + y_1 = z_1,$$
$$x_2 + y_2 = z_2,$$
$$x_3 + y_3 = z_3,$$
$$x_1 + x_2 = x_3,$$
$$y_1 + y_2 = y_3,$$
$$z_1 + z_2 = z_3.$$

Actually, there are only five *independent* equations; one of them (e.g. the last) can be derived from the others. (This corresponds to the fact that $z_3$ is the sum of both its column and its row, which overdetermines it.)

There are nine unknowns and five equations; that means we need four extra pieces of information (preferably, four of the values of the unknowns). That's exactly what the problem gives you; see the first table above. No matter which four variables we give you, you can solve for the others by a similar process.
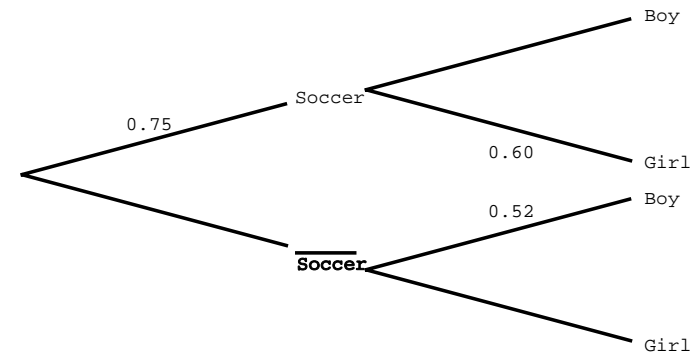
This problem can also be solved using a Venn diagram.  □

*Problem 2.* (15 points) In a certain school, 75% of the students play in the soccer league. Of the students who play in the league, 60% are girls. Of the students who do not play in the league, 52% are boys.
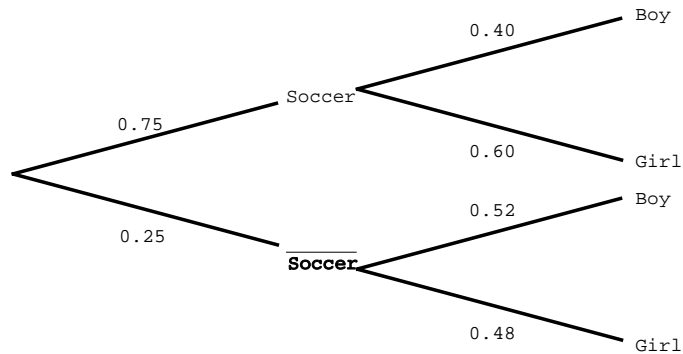
(a) What is the probability a randomly chosen student is a boy?
(b) What is the probability that a random chosen student from the school is a girl who plays in the soccer league?
(c) Mary's daughter attends this school. Knowing nothing else about her, what is the probability she plays in the soccer league?

*Solution.* This can also be solved as a table problem—although it seems we only have three pieces of information, not four, that's misleading; since these are percentages we know implicitly that the lower-right box of the table (the total of all the table entries) has to be 100%. However, most students will probably be more comfortable with the *tree* solution, which is the one we'll give.
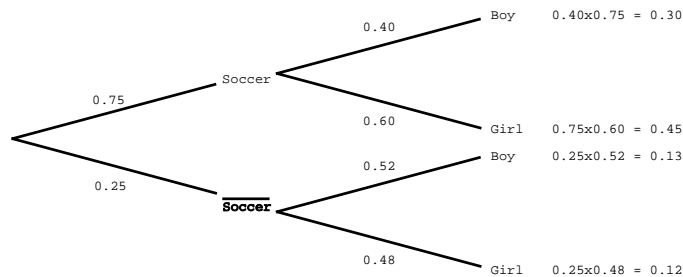
We'll start off with a tree, classifying students as to whether they play soccer ("Soccer") or don't ("$\overline{\text{Soccer}}$"); then as to whether they are boys or girls. The information we have can be summarized in our first tree:



Of course, if 75% of the students play soccer, then 25% don't; similarly, that which is not a girl is a boy, and that which is not a boy is a girl; with such commonplaces we can fill in some more of the tree:

Finally, we complete the tree by multiplying the probabilities of branches:



*Now* we're ready to answer the questions.

(a) The probability a student is a boy is $0.30 + 0.13 = 0.43$.
(b) The probability a randomly chosen student is a girl who plays in the soccer league is 0.45.
(c) We're asking for $P(S|G)$, the probability that a student plays soccer given that she's a girl. ("Mary" has nothing to do with it; it's the gender of her daughter which is the important datum.) Now

$$P(\text{Soccer}|G) = \frac{P(\text{Soccer and } G)}{P(G)} = \frac{0.45}{0.57} = 0.789.$$

□

*Problem* 3. (15 points) The Hiawatha City Council has fifteen members, nine of whom are pro-environment and six of who are pro-development. They wish to form a committee, consisting of five members of the city council, to consider new zoning legislation.

(a) In how many ways can the city council form such a committee?

(b) In how many ways can the city council form such a committee so that a majority of the committee is pro-environment?
(c) If the council selects the committee members randomly, what is the probability that exactly two of the members are pro-development?

*Solution.* Can you spell h-y-p-e-r-g-e-o-m-e-t-r-i-c? You can, of course, memorize (or bring in on your handwritten sheet) the relevant formulas–but why bother? All you have to remember is how many ways you can choose $r$ objects from $n$, and keep in mind the multiplication principle.

(a) There are 15 members of the council, and they want to choose 5 of themselves. The number of ways of doing this is

$$\binom{15}{5} = \frac{15 \times 14 \times 13 \times 12 \times 11}{1 \times 2 \times 3 \times 4 \times 5}$$
$$= 3003.$$

This is a product of five consecutive numbers (15 on down) divided by a product of five consecutive numbers (1 on up). The binomial coefficient $\binom{15}{5}$ is even *pronounced* "15 choose 5".

Alternatively, you can do it *ab initio,* from the multiplication principle: there are 15 ways of choosing the first committee member; leaving 14 ways to choose the second; leaving 13 ways to choose the third; leaving 12 ways to choose the fourth; leaving 11 ways to choose the fifth. That's a total of $15 \times 14 \times 13 \times 12 \times 11 = 360360$ ways of choosing the committee members, but it *counts order.* Since the hallmark of a committee is that it's just a *set* of people, without any order, we have to divide this by $5! = 120$, since there are 5! ways of permuting the order of any particular group of five people. And $360360/120 = 3003$.

(b) In order for a majority of the committee to be pro-environment, either 3 or 4 or 5 of the members must be pro-environment; in which case, 2 or 1 or 0 of the members will be pro-development. That is, we may have

   3 pro-environment, 2 pro-development
   4 pro-environment, 1 pro-development
   5 pro-environment, 0 pro-development.

Let's count the ways to do this. For example, how can we select 3 pro-environment and 2 pro-development members? Well, by selecting 3 of the 9 pro-environment members of the city council, and then selecting 2 of the 6 pro-development members. There are $\binom{9}{3} = 84$ ways of choosing the pro-environment members, and $\binom{6}{2} = 15$ ways of choosing the pro-development members; by the multiplication principle, there are $84 \times 15 = 1260$ ways of doing both. We summarize the number of ways in:

3 pro-environment, 2 pro-development: $\binom{9}{3}\binom{6}{2} = 84 \times 15 = 1260$

4 pro-environment, 1 pro-development: $\binom{9}{4}\binom{6}{1} = 126 \times 6 = 756$

5 pro-environment, 0 pro-development: $\binom{9}{5}\binom{6}{0} = 126 \times 1 = 126$

So there are a total of $1260 + 756 + 126 = 2142$ ways of choosing the committee to have a majority which is pro-environment.

(c) Choosing exactly *two* of the members to be pro-development can be done in exactly $1260$ ways (see part (b) above, under 3 pro-environment, 2 pro-development). There were $3003$ ways of choosing the committee altogether. Therefore, if the committee is picked at random, the probability that exactly two members will be pro-development is $1260/3003 = 60/143 = 0.41958$.

□

*Problem* 4. (25 points) Slick Motor Tire Inc. uses a tire mold which occasionally, and randomly, produces out-of-round tires. In fact 20% of the tires which it produces are out-of-round. Suppose Slick takes a random sample of 10 tires produced by this tire mold.

  (a) Find the expected number of out-of-round tires among these 10.
  (b) Find the probability that exactly two of these 10 tires are out-of-round.
  (c) Find the probability that two or more of these 10 tires are out-of-round.
  (d) Suppose now that the Fraud Motor Company buys 120 of the tires produced by this tire mold. Find the probability that between 20 and 30 (inclusive) of these tires are out-of-round.

*Solution.* This is a binomial probability problem, with $p = 0.2$. In items (a)-(c) we are dealing with a small population, 10, and we expect you to compute the binomial probabilities exactly. The point of part (d) is that the exact calculation, by hand, is onerous, so you are expected to approximate the binomial by a normal (and make the continuity correction).

  (a) Let $X$ be the total number of out-of-round tires among the ten. Then $E(X) = np = 10 \times 0.2 = 2$.
  (b) The probability that exactly two are out-of-round is

$$\binom{10}{2}0.2^2 0.8^8 = 45 \times 0.04 \times 0.167772 = 0.30199.$$

  (c) The probability that two *or more* are out-of-round is

$$\sum_{x=2}^{10}\binom{10}{x}0.2^x 0.8^{10-x},$$

which is a long calculation. Let's compute the *complementary* probability instead, that the number of out-of-round tires is 0 or 1. This is

$$\binom{10}{0}0.2^0 0.8^{10} + \binom{10}{1}0.2^1 0.8^9$$
$$= 1 \times 1 \times 0.107374 + 10 \times 0.2 \times 0.134218$$
$$= 0.107374 + 0.268436$$
$$= 0.37581.$$

This is the *complementary* probability, remember; the answer to the problem is therefore $1 - 0.37581 = 0.62419$.

  (d) Let $X$ be the number of the 120 tires which are out-of-round. We are asking for $P(20 \le X \le 30)$. We will use the normal approximation, *which means we must compute*

$$P(19.5 \le X \le 30.5)$$

instead. To do this we must compute

$$E(X) = np = 120 \times 0.2 = 24$$

and

$$\sigma_X = \sqrt{np(1-p)} = \sqrt{120 \times 0.2 \times 0.8} = 4.38178.$$

Then parallelling the construction of the standard normal variable $Z = (X - 24)/4.38178$,

$$P(19.5 \le X \le 30.5) = P(-4.5 \le X - 24 \le 6.5)$$
$$= P\left(\frac{-4.5}{4.38178} \le \frac{X - 24}{4.38178} \le \frac{6.5}{4.38178}\right)$$
$$= P(-1.02698 \le Z \le 1.48342)$$
$$= P(0 \le Z \le 1.02698) + P(0 \le Z \le 1.48342)$$
$$= 0.347785 + 0.431019$$
$$= 0.778804.$$

This uses an exact calculation of the normal distribution, but all you had available was the tables attached to the exam; thus an acceptable answer was

$$P(0 \le Z \le 1.03) + P(0 \le Z \le 1.48) = 0.3485 + 0.4306$$
$$= 0.7791.$$

(You could do a little bit better by linear interpolation, but this wasn't required; 0.7791 is fine.)

The *exact* answer is

$$\sum_{x=20}^{30} \binom{120}{x} 0.2^x 0.8^{120-x} = 0.776164,$$

available only if you had a laptop running MINITAB or Mathematica, or a premium calculator; but we see that 0.7791 was excellently close. On the other hand, if you *forgot* to make the continuity correction, and computed

$$P(20 \le X \le 30) = P(-4 \le X - 24 \le 6)$$

$$= P\left(\frac{-4}{4.38178} \le \frac{X-24}{4.38178} \le \frac{6}{4.38178}\right)$$

$$= P(-0.912871 \le Z \le 1.36931)$$

$$= P(0 \le Z \le 0.912871) + P(0 \le Z \le 1.36931)$$

limiting the precision for table look-up to

$$P(0 \le Z \le 0.91) + P(0 \le Z \le 1.37) = 0.3186 + 0.4147 = 0.7333,$$

your answer was farther off–and you lost points.

$\square$

*Problem* 5. (20 points) John regularly plays a target shooting competition in his yard. He is not very good at the game and his score has the following probability distribution

| Score | 0 | 1 | 2 | 3 |
|-------|-----|-----|-----|-----|
| Prob | 0.1 | 0.4 | 0.3 | 0.2 |

(a) Find the expected value and standard deviation of John's score.
(b) Suppose he plays the game twice. Find the probability that his total is at least 5.
(c) Suppose he plays the game 50 times. Find the probability that his total is at least 75.

*Solution.* The last part is designed to see whether you understand the Central Limit Theorem.

(a) The expected value is

$$0 \times 0.1 + 1 \times 0.4 + 2 \times 0.3 + 3 \times 0.2 = 1.6$$

The variance is

$$(0-1.6)^2 \times 0.1 + (1-1.6)^2 \times 0.4 + (2-1.6)^2 \times 0.3 + (3-1.6)^2 \times 0.2 = 0.84,$$

hence the standard deviation is $\sqrt{\text{Var}} = \sqrt{0.84} = 0.916515$.

(b) If he plays the game twice, the only outcomes where he scores five or more are (2, 3), (3, 2) and (3, 3). Assuming the plays are independent, the total probability of these is therefore

$$0.3 \times 0.2 + 0.2 \times 0.3 + 0.2 \times 0.2 = 0.16.$$

You can also draw a tree diagram for two plays, but there will be 16 branches in the second stage, and most of these will be wasted effort (since only three of them lead to a score of five or more).

(c) Ahhh, this is the heart of the problem! Let $X$ denote the number of points in 50 plays of the game. Assuming the plays are independent, we can write

$$X = \sum_{i=1}^{50} X_i,$$

where $X_i$ is the number of points scored the $i$-th time the game is played. Since the $X_i$ are independent identically distributed (i.i.d.) random variables, therefore

$$E(X) = 50 \times E(X_1) = 50 \times 1.6 = 80,$$

$$\sigma_X = \sqrt{50}\sigma_{X_1} = \sqrt{50} \times 0.916515 = 6.48074.$$

*By the Central Limit Theorem,* $X$ *is approximately a normally distributed random variable;* thus we are justified in setting
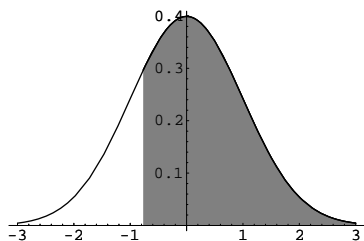
$$Z = \frac{X - 80}{6.48074},$$

and we compute

$$P(X \ge 75) = P(X - 80 \ge -5)$$

$$= P\left(\frac{X-80}{6.48074} \ge \frac{-5}{6.48074}\right)$$

$$= P(Z \ge -0.771517).$$

If you have a good calculator (one that can compute the cumulative distribution function for the standard normal variable), this works out to be 0.2202. If you must use the tables attached to the exam, then use

$$P(Z \ge -0.77) = 0.5 + P(0 \le Z \le 0.77)$$

$$= 0.5 + 0.2794$$

$$= 0.7794.$$

The explanation of why 0.5 is added to $P(0 \le Z \le 0.77)$ is best understood by referring to the following graph; we want the area of the shaded region:

The Central Limit Theorem says that $X$ is approximately normally distributed, but it doesn't say anything about how *good* the approximation is (at least, the version you've learned doesn't). The *exact* answer for this problem is $0.8016198\ldots$, so the answer yielded by the Central Limit Theorem is only about 3% off. How did I get the exact answer? Well, it's not something we covered in class. Here's the trick.

Form the polynomial

$$g(x) = 0.1 + 0.4x + 0.3x^2 + 0.2x^3,$$

obtained by using the *probabilities* as the coefficients of $x$ to the corresponding *value* of the random variable. (Refer back to John's scoring table to see how this is done.) Now square $g(x)$:

$$g(x)^2 = 0.01 + 0.08x + 0.22x^2 + 0.28x^3$$
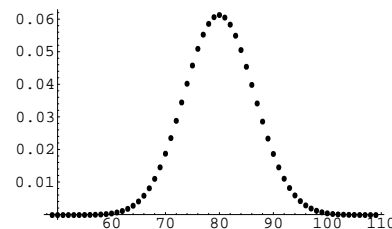$$+ 0.25x^4 + 0.12x^5 + 0.04x^6.$$

If you examine the tree diagram for the sum of two random variables $X_1$ and $X_2$ which both have John's scoring probability distribution, you obtain the table

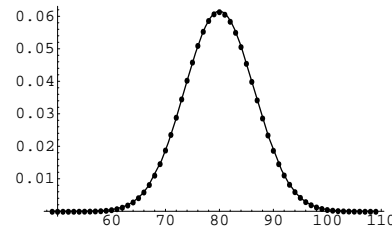| $X_1 + X_2$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|
| Prob | 0.01 | 0.08 | 0.22 | 0.28 | 0.25 | 0.12 | 0.04 |

The entries in the table are exactly the coefficients of $g(x)^2$!

This trick works more generally; in particular, if you want to find the probability distribution for $X_1 + \cdots + X_{50}$, where all the $X_i$ have John's scoring distribution, it suffices to compute the coefficients of $g(x)^{50}$. This is a polynomial of degree 150, and you really need a computer to do this, but the computer finds it quickly and easily.

Here's the actual probability density function for $X = X_1 + \cdots + X_{50}$, as a scatterplot:



Here's the same scatterplot, superimposed on the graph of the normal distribution with mean 80 and standard deviation 6.48074:
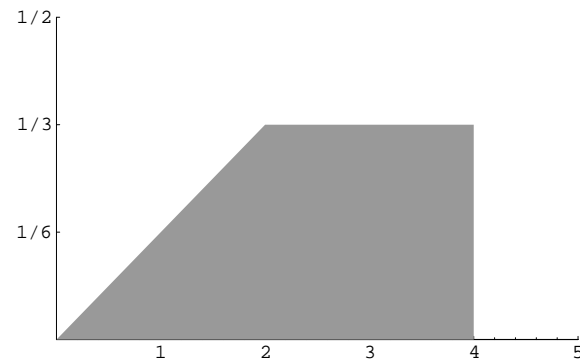


Yes, they're close!

When you add up the coefficients of $x^{75}, x^{76}, \ldots, x^{150}$ in the expansion of $g(x)^{50}$ you get 0.8016198. And since those coefficients were exactly the probability that John scored 75, 76, $\ldots$, or 150 points, 0.8016198 is exactly the probability that John scores 75 or more points.
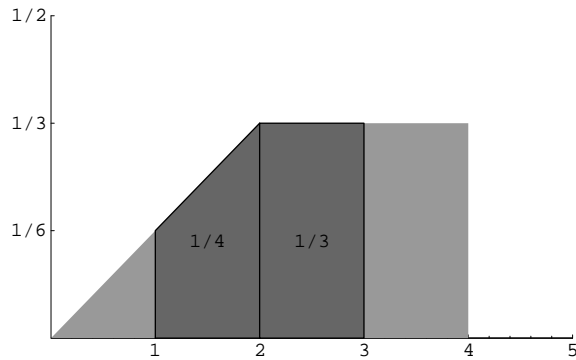
□

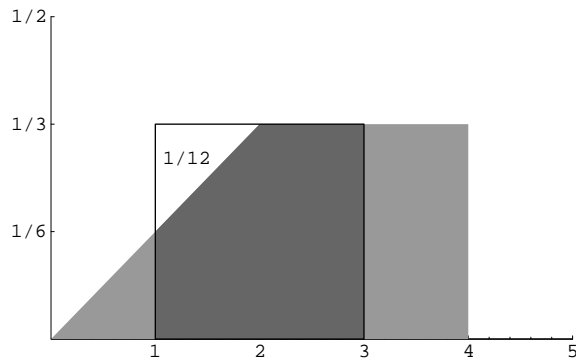*Problem 6.* (10 points) A random variable $X$ has probability density function shown below:



(a) Find the probability that $X$ is between 1 and 3.

(b) Find a value of $c$ such that $P(X \le c) = 0.95$.

*Solution.* First let's check to be sure this really *is* a probability density. The shaded area should have area 1; the triangular part has height 1/3 and base 2, so its area is $1/2 \times 1/3 \times 2 = 1/3$; the rectangle has height 1/3 and base 2, so its area is $1/3 \times 2 = 2/3$; indeed, the sum of the areas of the triangle and the rectangle is 1.
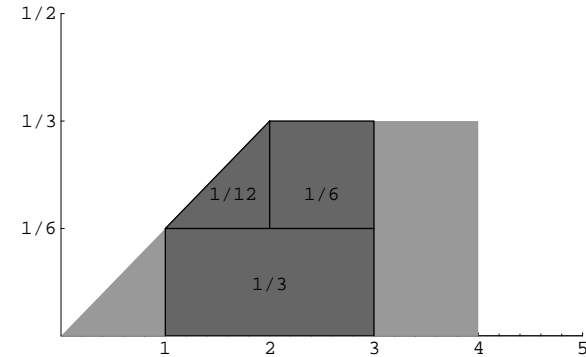
(a) You can calculate the area of the shaded region between $X = 1$ and $X = 3$ in any of several ways; for example, by dissecting it into a trapezoid with height 1 and sides 1/6 and 1/3 (turn your head sideways to see this), plus a rectangle with base 1 and height 1/3:



Or, you can compute the area of the rectangle from $x = 1$ to $x = 3$ and height 1/3, and subtract off the area of the small white triangle between $x = 1$ and $x = 2$ (which has area 1/12, since it has base 1 and height 1/6):



Or, you can think of it as the union of two rectangles and a triangle:

You *could* even compute it as

$$\int_1^2 \frac{1}{6}x\,dx + \int_2^3 \frac{1}{3}\,dx = \frac{1}{12}x^2 \Big|_{x=1}^3 + \frac{1}{3}x \Big|_{x=2}^3$$

However you dissect it, the area is always

$$\frac{7}{12} = 0.583333.$$

(a) The value of $P(X \le 2)$ is the area of the triangular shaded region, which is $1/3 = 0.3333$ and is far less than the desired 0.95; thus $c$ will lie in the interval $[2, 4]$. Since $P(X \le c) = 0.95$, we must have $P(X > c) = 0.05$, the complementary probability; the corresponding region is a rectangle whose height is 1/3 and whose base is $4 - c$. Thus we must have

$$\frac{1}{3}(4 - c) = 0.05,$$

a linear equation in $c$ which we easily solve to get $c = 3.85$. Or just use common sense: a rectangle of height 1/3 which has area 0.05 must have base $0.05/(1/3) = 0.15$; subtract 0.15 from 4.00 to get 3.85, the required value of $c$.

□

*Problem 7.* (15 points) The weights of packages carried by *United Envelope Express* courier service have a normal distribution with a mean of 17.2 ounces and a standard deviation of 6.3 ounces.

(a) It is proposed that packages weighing over 24 ounces will be subjected to a surcharge. What proportion of packages will be affected by the surcharge?

(b) A random sample of 20 packages is taken. Find the probability that their average weight is between 15 and 18 ounces.

*Solution.* Let $X$ be the random variable representing the weight of a randomly chosen package. Then

$$Z = \frac{X - 17.2}{6.3}$$

has the standard normal distribution.

(a) We are asked for $P(X \geq 24)$, and we compute it by parallelling the construction of $Z$ from $X$:

$$P(X \geq 24) = P(X - 17.2 \geq 6.8)$$

$$= P\left(\frac{X - 17.2}{6.3} \geq \frac{6.8}{6.3}\right)$$

$$= P(Z \geq 1.07937)$$

$$= 0.140211$$

if you have a calculator which can compute the cumulative normal distribution accurately; otherwise, if you use the tables, use

$$P(Z \geq 1.08) = 0.5 - 0.3599 = 0.1401.$$

Just over 14% of the packages will be subjected to the surcharge.

(b) Let $\overline{X}$ be the average weight of 20 packages. Then $\overline{X}$ is normally distributed with mean $E(\overline{X}) = 17.2$ and standard deviation

$$\sigma_{\overline{X}} = \sigma_X/\sqrt{20} = 6.3/\sqrt{20} = 1.40872.$$

Therefore

$$Z = \frac{\overline{X} - 17.2}{1.40872}$$

has the standard normal distribution, and we compute

$$P(15 \leq \overline{X} \leq 18) = P(-2.2 \leq \overline{X} - 17.2 \leq 0.8)$$

$$= P\left(\frac{-2.2}{1.40872} \leq \frac{\overline{X} - 17.2}{1.40872} \leq \frac{0.8}{1.40872}\right)$$

$$= P(-1.5617 \leq Z \leq 0.567891)$$

$$= 0.655766,$$

at least, if you have a calculator good enough to compute the CDF. If you have to do table lookup, then

$$P(-1.56 \leq Z \leq 0.57) = P(0 \leq Z \leq 1.56) + P(0 \leq Z \leq 0.57)$$

$$= 0.4406 + 0.2157 = 0.6563$$

will do fine. (The original data were only given to three significant figures, so going for six significant figures in our answer would be silly.)

$\square$

*Problem 8.* (20 points) You are constructing a marketing research study to determine what proportion of consumers between the ages of 18 and 35 own DVD players. You have randomly sampled 60 individuals in that age range, and you have determined that 11 of them own DVD players.

(a) Find a point estimate for the proportion of consumers who own DVD players.
(b) Find the 90% confidence interval for the proportion.
(c) Suppose that your contract for the study requires that you estimate the proportion within plus or minus 0.05 with 90% confidence. How many additional subjects would you need to sample to obtain that level of precision?

*Solution.*

(a) The point estimate for $p$ is

$$\hat{p} = 11/60 = 0.1833.$$

(b) Since this is a binomial probability problem, we will use the normal distribution to approximate it. The 90% cutoff for $Z$ is

$$Z_{0.05} = 1.645.$$

The confidence interval has endpoints

$$\hat{p} \pm Z_{0.05}\sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

$$= 0.1833 \pm 1.645\sqrt{\frac{0.1833 \times 0.8167}{60}}$$

$$= 0.1833 \pm 0.0822.$$

So $L = 0.1011$ and $R = 0.2655$.

(c) The formula for $n$ is

$$n = \frac{4Z_{\alpha/2}^2\hat{p}(1 - \hat{p})}{w^2}$$

$$= \frac{4 \times 1.645^2 \times 0.1833 \times 0.8167}{0.1^2}$$

$$= 162.038.$$

Since only an integer number of consumers can be questioned, we must round up to $n = 163$. The problem asks for the *additional* number of consumers who must be questioned (beyond the original 60), hence the correct answer is 103.

$\square$

*Problem* 9. (20 points) Assume that the weight of the cereal in boxes of *Loopy Froots* breakfast cereal is a random variable that is normally distributed. A consumer advocate group has randomly sampled 6 boxes of *Loopy Froots* from supermarket shelves and weighed the contents. The table below shows the data that they collected.

| Weight in grams | 535 | 540 | 565 | 575 | 535 | 550 |
|---|---|---|---|---|---|---|

(a) Find point estimates for the population mean and standard deviation.
(b) Find a 95% confidence interval for the population mean.
(c) Suppose that the mean weight of a box of *Loopy Froots* is supposed to be 570 grams. Based on your result in part (b), does the consumer advocate group have reason to believe that the actual mean weight is different from 570 grams? Please explain.

*Solution.*

(a) We compute the mean $\overline{X}$ to be 550 and the sample standard deviation $s$ to be $2\sqrt{70} = 16.7332$. The latter is computed by

$$s^2 = \frac{1}{5}\big((535-550)^2 + (540-550)^2 + (565-550)^2$$
$$+ (575-550)^2 + (535-550)^2 + (550-550)^2\big)$$

Notice the division by 5, not 6!

(b) We use the Student $t$-statistic with 5 degrees of freedom (since there are $n = 6$ measurements). We find from the table that the cutoff is $t_{0.025,5} = 2.571$. The confidence interval therefore has endpoints

$$\overline{X} \pm t_{\alpha/2,n-1}\frac{s}{\sqrt{n}}$$
$$= 550 \pm 2.571\frac{16.7332}{\sqrt{6}}$$
$$= 550 \pm 17.5633$$

leading to a 95% confidence interval with left-hand endpoint $L = 532.437$ and right-hand endpoint $R = 567.5633$.

(c) *This item is the crux of the problem!* What's the point of doing all the technical computations if you *don't reach the right conclusion?*

*Yes,* the consumer advocate group has reason to believe that the actual mean weight is different from 570 grams. (The answer is *not* "No, the consumer advocate group has reason to believe that the actual mean weight is different from 570 grams." That defies the meanings of the words yes and no.)

The reason is that we're 95% sure that the true mean is somewhere between 532.4 and 567.6 grams, and 570 grams isn't between these. We're *at least* 95% sure that the true mean is less than 570 grams.

Alternatively, we compute the $t$ statistic as

$$t = \frac{\overline{X} - 570}{s/\sqrt{n}}$$
$$= \frac{550 - 570}{16.7332/\sqrt{6}}$$
$$= -2.9277.$$

Since $t < -t_{0.025,5} = -2.571$, we should *reject* any suggestion that the mean is 570. The packages are significantly underweight, and the feds are coming . . .

The small sample size has *nothing to do with it*. Small sample sizes are taken account of in the larger values of the $t$-cutoffs. Of course, if the true mean had been less than 570 but much closer to 570, we would have needed a larger sample; but in this problem the difference was already so extreme that a larger sample wasn't necessary. Think of it this way: if we buy half-a-dozen packages of elephant meat, and we're told each package weighs a thousand pounds; but when we weigh them, they all weigh less than a hundred pounds, we don't bother buying any more sample packages–we call the consumer protection agency.

□

*Problem* 10. (20 points) When the machine at Golden Valley Brewing Co. is properly adjusted the volume of beer it puts in each bottle is normally distributed with a mean of 12.1 oz. When the machine gets out of adjustment the mean may be higher or lower. Every day the machine is tested to determine if it is properly adjusted. This is done by taking a sample of fifteen beer bottles and measuring the amount of beer in each bottle. If the test indicates, at the 1% significance level, that the mean has changed from 12.1 oz, then production is halted to allow the machine to be readjusted.

(a) Formulate the null and alternative hypotheses associated with the daily tests.
(b) Choose an appropriate test statistic and find the rejection rule for the daily tests.
(c) On Monday, the sample average is 12.06 and the sample standard deviation is 0.05 oz. Should production be halted to allow the machine to be readjusted? Justify your decision.
(d) On Tuesday, following that day's test, production is halted to allow the machine to be readjusted. However, on close inspection it is found that

the machine does not need readjustment. Determine which of the two statements below best describes this situation.

(i) The test must have been conducted improperly (for example the amount of beer in one of the sample bottles may have been mis-measured) and that resulted in the error.

(ii) The nature of the testing procedure is that errors of this type will happen occasionally even if the test is conducted perfectly.

*Solution.*

(a) Null hypothesis $H_0$: $\mu = 12.1$, where $\mu$ is the average amount of beer (in ounces) going into the bottles. Alternative hypothesis $H_a$: $\mu \neq 12.1$. We will be led to reject the null hypothesis if the machine is significantly underfilling *or* overfilling the bottles.

(b) The test statistic will be the Student $t$-distribution, computed by

$$t = \frac{\overline{X} - 12.1}{s/\sqrt{n}},$$

where $\overline{X}$ is the average in the 15 sample bottles, $s$ is the sample standard deviation for those 15 bottles, and $n = 15$. We will *reject* the null hypothesis (at the 1% significance level) if $|t| > t_{0.005,14} = 2.977$.

(c) Let's compute the $t$-statistic with these particular values:

$$t = \frac{\overline{X} - 12.1}{s/\sqrt{n}}$$
$$= \frac{12.04 - 12.1}{0.05/\sqrt{15}}$$
$$= -3.09839.$$

Since $t < -2.977$, we *should* reject the null hypothesis.

(d) Statement (ii) is correct. Statistical deductions are just that: statistical; at the 1% level of significance, we expect to be wrong 1% of the time, just because we accidentally got extreme data. "My momma told me there'd be days like this."

□

*Problem* 11. (20 points) Obstructive sleep apnea is a disorder that causes a person to stop breathing and awaken briefly during a sleep cycle. About 25% of the population suffers from this disorder. Researchers suspect that truck drivers may be more prone to the disorder than the rest of the population. To check this hypothesis they select a random sample of 59 truck drivers and assess how many of them suffer from the disorder.

(a) Formulate the null and alternative hypotheses that the researchers use.

(b) It is found that 22 of the truck drivers selected suffer from the disorder. Calculate the value of a suitable test statistic and find the P-value.

(c) At which of the following levels of significance can the null hypothesis be rejected? Circle all that apply:

    0.2     0.1     0.05     0.02     0.01

*Solution.* Since I suffer from sleep apnea, I found this to be a fascinating problem. The data are real, incidentally; apparently truck drivers *do* tend to suffer from sleep apnea more than the general population. It's not a comforting thought that the driver of that 18-wheeler approaching me on the highway may not only be asleep but also not breathing . . .

(a) The null hypothesis is $H_0 : p = 0.25$, where $p$ is the proportion of truck drivers who suffer from sleep apnea; the alternative hypothesis, since the researchers *expect* that drivers are more prone to the illness, is $p > 0.25$. (One could also take as the null hypothesis $H_0 : p \leq 0.25$.)

(b) The results of the test find $\hat{p} = 22/59 = 0.372881$, and we form the test statistic

$$Z = \frac{\hat{p} - 0.25}{\sqrt{0.25 \times 0.75/n}} = 2.17977.$$

The P-value is $P(Z \geq 2.17977) = 0.0146373$. (If you use the tables, then the P-value is $P(Z \geq 2.18) = 0.0146$.)

(c) We can reject at any significance level which *exceeds* 0.0146; e.g. at the 2% level, the 5% level, the 10% level and the 20% level. The only level among those listed at which we *cannot* reject the null hypothesis is the 1% level.

□