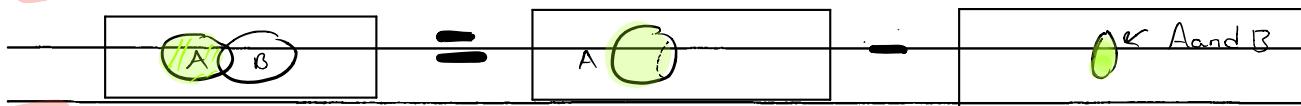


Rules (Contd.)

8) $\Pr\{A \text{ and } B^c\} = \Pr\{A\} - \Pr\{A \text{ and } B\}$



9) $\Pr\{A^c | B\} = 1 - \Pr\{A | B\}$ or

$$\Pr\{B^c | A\} = 1 - \Pr\{B | A\}$$

BUT $\Pr\{A | B^c\} \neq 1 - \Pr\{A | B\}$ and $\Pr\{B | A^c\} \neq 1 - \Pr\{B | A\}$

For example, if we know that out of all women, 25% wear glasses, we also know out of all women, 75% do not wear glasses ($\Pr\{G^c | W\} = 1 - \Pr\{G | W\}$). However, we can not say that 75% of men wear glasses ($\Pr\{G | W^c\}$).

Ex: Suppose the prob. of someone swimming regularly is 0.17, and the prob. of someone swimming and getting an ear infection within a year is 0.09. Find the following:

a) The prob. that someone swims and does not get an ear infection (within a year).

Let $S = \text{Swim}$, $I = \text{Infection}$.

$$\Pr\{S \text{ and } I^c\} = \Pr\{S\} - \Pr\{S \text{ and } I\} = 0.17 - 0.09 = 0.08$$

b) If someone swims, the prob. they get an infection.

$$\Pr\{I | S\} = \Pr\{I \text{ and } S\} / \Pr\{S\} = 0.09 / 0.17 = 0.529$$

c) The prob. someone does not get an infection, if they swim.

$$\Pr\{I^c | S\} = 1 - \Pr\{I | S\} = 1 - 0.529 = 0.471$$

Definition: Two events A and B are independent if and only if either of these conditions holds true:

i) $\Pr\{A \text{ and } B\} = \Pr\{A\} \Pr\{B\}$

ii) $\Pr\{A | B\} = \Pr\{A\}$

Independent means that the prob. of one event does not effect the prob. of another.

The equations (i) and (ii) can only be used if you know two events are independent. Generally, we do not assume they are.

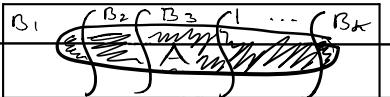
Last rule:

10) For any two events A and B,

$$\Pr\{A\} = \Pr\{A \text{ and } B\} + \Pr\{A \text{ and } B^c\}$$

For any events A and $B_1, B_2, B_3, B_4, \dots, B_k$,

$$\Pr\{A\} = \Pr\{A \text{ and } B_1\} + \Pr\{A \text{ and } B_2\} + \dots + \Pr\{A \text{ and } B_k\}$$

 We are just adding up all parts of B that overlap with A, to get A.

This rule is sometimes called the "rule of total prob".

Ex: The probability of a certain disease is 0.05. If a person has the disease, the prob of testing positive is 0.99. If they do not have the disease, the prob of testing positive is 0.02.

a) Write down in probability notation the values you have been given.

Let + = positive, D = disease.

$$\text{Given: } \Pr\{D\} = 0.05, \Pr\{+ | D\} = 0.99, \Pr\{+ | D^c\} = 0.02.$$

$$\text{We can also get: } \Pr\{D^c\} = 0.95, \Pr\{- | D\} = 0.01, \Pr\{- | D^c\} = 0.98$$

b) Find the prob that someone tested positive.

$$\Pr\{+\} = \Pr\{+ \text{ and } D\} + \Pr\{+ \text{ and } D^c\} \quad (\text{using rule 10})$$

$$= \Pr\{+ | D\} \Pr\{D\} + \Pr\{+ | D^c\} \Pr\{D^c\} \quad (\text{using rule 7 twice})$$

$$= (0.99)(0.05) + (0.02)(0.95) = 0.0685$$

* Why is this prob. larger than $\Pr\{D\}$?

c) What is the prob that someone who tested + actually has the disease?

$$P_{\text{Pr}} \{\xi + \text{D}\} = \frac{P_{\text{Pr}} \{\xi + \text{D}\} P_{\text{Pr}} \{\xi\}}{P_{\text{Pr}} \{\xi\}} = \frac{0.99(0.05)}{0.0685} = 0.7726$$

I.e., if you have tested positive, there is a 77.26% chance you have the disease.

d) What is the probability that someone tests positive, or has the disease, or both?

$$P_{\text{Pr}} \{\xi + \text{or D}\} = P_{\text{Pr}} \{\xi + \text{D}\} + P_{\text{Pr}} \{\xi\} - P_{\text{Pr}} \{\xi + \text{and D}\} = 0.05 + 0.0685 - 0.99(0.05) \\ = 0.069$$

e) Are testing positive and having the disease independent?

To check this, we would see if

$$\text{i) } P_{\text{Pr}} \{\xi + \text{and D}\} = P_{\text{Pr}} \{\xi + \text{D}\} P_{\text{Pr}} \{\xi\} \quad \text{or ii) } P_{\text{Pr}} \{\xi + \text{D}\} = P_{\text{Pr}} \{\xi\}$$

(you can check any other combinations as well)

$$\text{For i: } P_{\text{Pr}} \{\xi + \text{D}\} P_{\text{Pr}} \{\xi\} = 0.0685(0.05) = 0.003425 \\ P_{\text{Pr}} \{\xi + \text{and D}\} = 0.99(0.05) = 0.0495$$

These are not equal, so they are not independent
(they are dependent).

For ii) $P_{\text{Pr}} \{\xi + \text{D}\} = 0.99$, $P_{\text{Pr}} \{\xi + \text{D}\} = 0.0685$, which are not equal.

(note; you only needed to check one of these).



Discrete Random Variables

Recall: A discrete r.v. is numeric with natural gaps.

Fact: All the probability rules we learned also apply to all r.v's.

Ex: Let $Y = \#$ on a die. Then, Y can equal 1, 2, 3, 4, 5, 6, so is a discrete r.v.

Notice all values of Y are mutually exclusive. This is true for all r.v's.

We can use our probability rules to calculate:

a) $P_{\text{Pr}} \{Y > 5\}$ (we roll a number > 5) = $P_{\text{Pr}} \{Y = 6\}$
(because of the gap) = $\frac{1}{6}$ (if it is a fair die)

$$\text{“A} \cap \text{“} = P\{\text{A}\}$$

$$\text{b) } P\{Y \geq 1\} = 1 - P\{Y < 1\} = 1 - P\{Y = 0\} = 1 - \frac{1}{6} = \frac{5}{6}$$

$$\text{c) } P\{Y \leq 1 \text{ or } Y \geq 5\} = P\{Y \leq 1\} + P\{Y \geq 5\} + 0 \quad (\text{mutually exclusive}) \\ = P\{Y = 1\} + P\{Y = 5\} + P\{Y = 6\} = \frac{3}{6}$$

In addition, for any discrete r.v we have,

Definition: The mean of a discrete r.v Y is (or expected value)

$$\mu_Y = \sum_{i=1}^k y_i P\{Y = y_i\}, \quad Y \text{ takes on values } y_1, y_2, \dots, y_k$$

Definition: The variance of a discrete r.v Y is

$$\sigma^2_Y = \sum_{i=1}^k (y_i - \mu_Y)^2 P\{Y = y_i\} = [\sum y_i^2 P\{Y = y_i\}] - (\mu_Y)^2$$

For the dice example, we have

$$\mu_Y = 1(\frac{1}{6}) + 2(\frac{1}{6}) + 3(\frac{1}{6}) + 4(\frac{1}{6}) + 5(\frac{1}{6}) + 6(\frac{1}{6}) = 3.5$$

$$\sigma^2_Y = [\sum y_i^2 P\{Y = y_i\}] - 3.5^2 \\ = [1^2(\frac{1}{6}) + 2^2(\frac{1}{6}) + 3^2(\frac{1}{6}) + 4^2(\frac{1}{6}) + 5^2(\frac{1}{6}) + 6^2(\frac{1}{6})] - 3.5^2 \\ = 15.166 - 12.25 = 2.916$$

Ex: In a year of observing patients, the distribution of time to "full recovery" after knee replacement was:

#months	6	12	18	24
Rel. Freq	0.20	0.30	0.35	0.15

You may assume this is population data

Let $Y = \# \text{ months to full recovery}$.

a) Find the average of Y .

$$\mu_Y = \sum y_i \Pr\{Y=y_i\} = 6(0.20) + 12(0.30) + 18(0.35) + 24(0.15) \\ = 14.7 \text{ months}$$

b) Find the variance of Y .

$$\sigma_Y^2 = \sum (y_i - \mu_Y)^2 \Pr\{Y=y_i\} = (16-14.7)^2(0.20) + (12-14.7)^2(0.30) + \\ (18-14.7)^2(0.35) + (24-14.7)^2(0.15) = 34.11$$

c) Interpret the standard deviation in terms of the problem.

$\sigma_Y = \sqrt{34.11} = 5.84$, a typical distance/deviation from the mean is 5.84 months.

* Linear Combinations of Random Variables

If X is a r.v. with mean μ_X , std dev σ_X , then

i) If $Y = X + a$ (we add a constant "a" to all X values)

$$\mu_Y = \mu_X + a, \quad \sigma_Y = \sigma_X$$

ii) If $Y = bX$ (we multiply every X value by a constant "b")

$$\mu_Y = b\mu_X, \quad \sigma_Y^2 = b^2 \sigma_X^2, \quad \sigma_Y = |b| \sigma_X$$

These give the last result,

iii) If $Y = a + bX$ (a linear combination)

$$\mu_Y = a + b\mu_X, \quad \sigma_Y^2 = b^2 \sigma_X^2, \quad \sigma_Y = |b| \sigma_X$$

Ex: If X has $\mu_X = 4$, $\sigma_X = 2$, and $Y = 2 - 4X$

Find μ_Y and σ_Y^2

$$a = 2, \quad b = -4, \quad \mu_Y = 2 + (-4)\mu_X = -12$$

$$\sigma_Y^2 = b^2 \sigma_X^2 = (-4)^2(2^2) = 64$$

* These properties apply to any r.v.



Binomial Random Variables

A very frequently used r.v. is one with the following properties

- 1) An experiment is conducted with exactly two possible outcomes, one of which is a "success", one a "failure".
- 2) The outcomes of all experiments/trials are independent.
- 3) The probability of a "success" for a single trial is denoted p , and p does not change.
- 4) The variable of interest is $Y = \# \text{ of successes in } n \text{ trials}$

Then, Y is a binomial r.v., which can take on values $0, 1, 2, 3, \dots, n$.

The parameters of a binomial distribution are $n = \# \text{ of trials}$, $p = P\{\text{"success"}\}$.

Ex: A coin flipped 5 times, where a "success" is a head.
 $n = 5$, $p = 0.50$

Ex: The number of children who are over 65 inches tall, out of 4. $n = 4$, $p = P\{\text{over 65 inches}\}$

Ex: A couple has 3 children, and each child has a 15% chance of having blue eyes.

a) Identify the distribution, and its parameters, if we are interested in the # of children with blue eyes.

$Y = \# \text{ of children with blue eyes out of 3}$, $n = 3$, $p = 0.15$. This is a binomial r.v., Y can = $0, 1, 2, 3$

b) Let $B_i = \text{event child } i \text{ has blue eyes}$. Find the probability that no children have blue eyes.

We want:

NOTE:
 $P_r\{\bar{B}_1 \text{ and } \bar{B}_2 \text{ and } \bar{B}_3\} = 1 - P_r\{B_1 \cup B_2 \cup B_3\}$

$$\Pr\{\bar{B}_1 \text{ and } \bar{B}_2 \text{ and } \bar{B}_3\} \stackrel{\checkmark \text{ by independence}}{=} \Pr\{\bar{B}_1\} \Pr\{\bar{B}_2\} \Pr\{\bar{B}_3\}, \\ = (1 - .15)(1 - .15)(1 - .15) = 0.614$$

c) Find The probability that exactly two children have blue eyes.

$$\Pr\{Y = 2\} = \Pr\{\{B_1 \text{ and } B_2 \text{ and } \bar{B}_3\} \text{ or } \{B_1 \text{ and } \bar{B}_2 \text{ and } B_3\} \text{ or } \{\bar{B}_1 \text{ and } B_2 \text{ and } B_3\}\}$$

$$\begin{aligned} (\text{Mutually Ex.}) &= \Pr\{\bar{B}_1 \text{ and } B_2 \text{ and } B_3\} + \Pr\{\bar{B}_1 \text{ and } \bar{B}_2 \text{ and } B_3\} + \Pr\{\bar{B}_1 \text{ and } B_2 \text{ and } \bar{B}_3\} \\ &= \Pr\{\bar{B}_1\} \Pr\{B_2\} \Pr\{\bar{B}_3\} + \Pr\{\bar{B}_1\} \Pr\{\bar{B}_2\} \Pr\{B_3\} + \Pr\{\bar{B}_1\} \Pr\{B_2\} \Pr\{\bar{B}_3\} \\ &= (0.15)(0.15)(0.85) + (0.15)(0.85)(0.15) + (0.85)(0.15)(0.15) = 3(0.15^2)(0.85) \\ &= 0.0574 \end{aligned}$$

I.e,

$$\Pr\{Y = 2\} = 3(0.15^2)(0.85)$$

$3 = \# \text{ of ways we can get 2 blue eyed children out of 3}$

$(0.15)^2 = \text{prob of getting 2 blue eyed children}$

$(0.85) = \text{prob of getting 1 non-blue eyed child}$

Notice we used that values of Y are mutually exclusive, while the "trials" (children) are independent.

If n were larger, it would be tedious to calculate the above. However, a binomial r.v. has simplified equations for calculating $\Pr\{Y = j\}$, where j is some number:

* Properties of binomial random variables

Let Y be a binomial r.v. with parameters n, p . Then

the following hold true:

$$i) \Pr\{Y = j\} = \binom{n}{j} p^j (1-p)^{n-j}, \text{ for } j = 0, 1, 2, \dots, n$$

= prob of j successes in n trials, where

$$\binom{n}{j} = \frac{n!}{j!(n-j)!}, \text{ where } n! = n(n-1)(n-2)\cdots(3)(2)(1)$$

Other notation for $\binom{n}{j}$ is nC_j . Note: $0! = 1$.

$$ii) \mu_Y = \text{expected / average } \# \text{ of successes in } n \text{ trials} \\ = np$$

$$iii) \sigma^2_Y = \text{variance of the } \# \text{ of successes in } n \text{ trials} \\ = np(1-p)$$

$$\text{Ex of } n! : 4! = 4(3)(2)(1) = 24, \quad \binom{4!}{2} = \frac{4!}{2!(4-2)!} = \frac{4 \cdot 3 \cdot 2 \cdot 1}{(2 \cdot 1)(2 \cdot 1)} = 3 \cdot 2 \cdot 1 = 6$$

Ex: A plant has a 20% chance of having a red flower bloom on any bud. If a plant has 6 buds, find the following:

a) The prob. of exactly one red flower.

Y is binomial, $n = 6$, $p = 0.20$

$$\Pr\{Y=1\} = \binom{6}{1} (0.20)^1 (1-0.20)^5 = \frac{6!}{1!(6-1)!} (0.20)(0.80)^5 = \frac{6 \cdot 5 \cdot 4 \cdot 3 \cdot 2 \cdot 1}{(1)5 \cdot 4 \cdot 3 \cdot 2 \cdot 1} (0.0655) \\ = 0.393$$

b) The prob. of at least two red flowers.

$$\Pr\{Y \geq 2\} = 1 - \Pr\{Y < 2\} = 1 - \Pr\{Y \leq 1\} = 1 - (\Pr\{Y=0\} + \Pr\{Y=1\}) \\ = 1 - (0.393) - \binom{6}{0} \cdot 0.20^0 (1-0.20)^6 = 1 - 0.393 - .80^6 \\ = 1 - 0.393 - 0.262 = 0.345$$

c) The mean and variance for the # of red flowers out of 6

$$\mu_Y = np = 6(0.20) = 1.2, \quad \sigma^2_Y = np(1-p) = 6(0.20)(0.80) = .96$$

Note: $\binom{n}{0} = 1$, $\binom{n}{1} = n$, $\binom{n}{n-1} = n$, for any n .

Do not round μ_Y , σ^2_Y .



Probabilities for Continuous Random Variables

With continuous data, we cannot use the same approach as we did with discrete, since there are an infinite number of values.

We instead find probabilities as areas under density curves.

You can think of a density curve as a smoothed histogram:



or



The curve is the function that creates the graph.

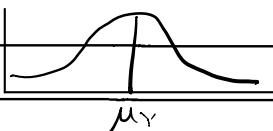
Density curves have the following properties:

- 1) The total area under the curve must be 1
- 2) Probabilities are areas under curves (found via calculus)
- 3) The area under the curve between points "a" and "b" gives the prob. the r.v. takes on values between "a" and "b".



Normal Random Variables

The most widely used continuous r.v. in statistics is the Normal / Gaussian r.v. T is perfectly symmetric about its mean, and is "bell shaped":



Some facts about a normal r.v.:

- 1) $\bar{y} = \tilde{y}$ (the mean equals the median)
- 2) $\Pr\{Y > \mu + a\} = \Pr\{Y < \mu - a\}$  (Symmetric)
- 3) If Y is distributed normal with mean μ_Y , std dev σ_Y , we denote it as $Y \sim N(\mu_Y, \sigma_Y^2)$



Standardizing data

For any random variable Y with mean μ_Y , std. dev. σ_Y , if we use the following linear transformation:

$$Y^* = \frac{(Y - \mu_Y)}{\sigma_Y}, \text{ then } \mu_{Y^*} = 0, \sigma_{Y^*} = 1.$$

To show this:

$$Y^* = \frac{Y}{\sigma_Y} - \frac{\mu_Y}{\sigma_Y}, \quad a = \frac{\mu_Y}{\sigma_Y}, \quad b = \frac{1}{\sigma_Y}, \quad \text{so } \mu_{Y^*} = \frac{\mu_Y}{\sigma_Y} - \frac{\mu_Y}{\sigma_Y} = 0$$

$$\sigma_{Y^*}^2 = b^2 \sigma_Y^2 = (\frac{1}{\sigma_Y})^2 \sigma_Y^2 = 1 \Rightarrow \sigma_{Y^*} = 1$$



Standard Normal Random Variable

When $Y \sim N(\mu_Y, \sigma_Y)$, the transformation $(Y - \mu_Y)/\sigma_Y$ is also normal, with mean 0, standard deviation 1. This is called a standard normal r.v., and typically denoted with Z .

I.e. If $Y \sim N(\mu_Y, \sigma_Y)$, and $Z = (Y - \mu_Y)/\sigma_Y$
Then $Z \sim N(0, 1)$

The "Z-table" in your book has calculated areas under curves for $Z \sim N(0, 1)$. The table gives $P_r\{Z < Z^*\}$, $Z^* = -3.49, -3.48, \dots, 3.48, 3.49$

$Z = (Y - \mu_Y)/\sigma_Y$ can be calculated on a single data point, at which point we call it a z-score.
z-score = $(Y - \mu_Y)/\sigma_Y$ = The # of std dev's Y is from μ_Y .

The table in the back gives the second decimal place as columns, the first two digits as rows, and probabilities in the middle.

		0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0		$P_r\{Z < 0.00\}$	$P_r\{Z < 0.04\}$								
0.1											
:											
1.0	0.8413										
1.1									$P_r\{Z < 1.15\}$		
:											
3.4										$P_r\{Z < 3.49\}$	
:											

Ex: $P_r\{Z < 1\}$ can be found directly in row 1.0, column 0.00 = 0.8413

$$P_r\{Z > 1\} = 1 - P_r\{Z < 1\} = 1 - 0.8413 = 0.1587 \quad (\text{Shaded area} = \text{Total area} - \text{Unshaded area})$$

$$P_r\{-1 < Z < 1\} = \text{Shaded area between } -1 \text{ and } 1 = \text{Total area} - \text{Area outside } (-1, 1) \quad \text{Too much}$$

$$= P_r\{Z < 1\} - P_r\{Z < -1\} = 0.8413 - 0.1587 = 0.6826$$

$$(\text{Notice } P_r\{Z > 1\} = P_r\{Z < -1\} \text{ by symmetry})$$



Percentiles for a Z , and $Z_{\alpha/2}$

The k^{th} percentile for a Z can be found by looking up $\frac{k}{100}$ (or the closest value to it) in the middle of the Z table, then seeing what Z -score that is by looking at what row and column it was in.

Ex: Find the $\sim 70^{\text{th}}$ percentile for a Z .

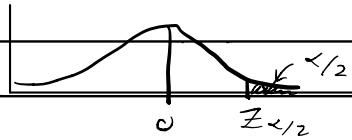
0.6985 is in column .02, row 0.50, so $P_r\{Z < .52\} = 0.6985$

0.7019 is in column .03, row 0.50, so $P_r\{Z < .53\} = 0.7019$

Since .6985 is closer to 0.70, the approximate 70^{th} percentile for Z is $Z^{(70)} = .52$.

Note: If two values in the middle are equi-distance from $\frac{k}{100}$, average the Z -scores.

$Z_{\alpha/2}$ is defined to be the value of Z with $\alpha/2$ in the upper tail;



So for ex, $Z_{.10}$ has 0.10 in the upper tail, 0.40 in the lower (it is the 90^{th} percentile).

We will use this notation later.

* Using the Z -table for all $Y \sim N(\mu_Y, \sigma_Y)$

When you have $Y \sim N(\mu_Y, \sigma_Y)$, and want any probability of Y , you can transform Y into $Z \sim N(0,1)$ then use the Z -table:

$$(i) P_r\{Y < a\} = P_r\left\{ \frac{(Y - \mu_Y)}{\sigma_Y} < \frac{(a - \mu_Y)}{\sigma_Y} \right\} = P_r\{Z < \frac{(a - \mu_Y)}{\sigma_Y}\}$$

$$(ii) P_r\{Y > a\} = 1 - P_r\{Y < a\} = 1 - P_r\left\{ \frac{(Y - \mu_Y)}{\sigma_Y} < \frac{(a - \mu_Y)}{\sigma_Y} \right\} = 1 - P_r\{Z < \frac{(a - \mu_Y)}{\sigma_Y}\}$$

$$(iii) P_r\{a < Y < b\} = P_r\{Y < b\} - P_r\{Y < a\} = P_r\{Z < \frac{(b - \mu_Y)}{\sigma_Y}\} - P_r\{Z < \frac{(a - \mu_Y)}{\sigma_Y}\}$$

(iv) To find a percentile for Y , first find the percentile for Z , then solve for Y :

$$Z^{(k)} = \frac{Y^{(k)} - \mu_Y}{\sigma_Y} \Rightarrow Y^{(k)} = \mu_Y + Z^{(k)}(\sigma_Y)$$

Ex: The distribution for the score on an exam is normal, with mean 60, std. dev 10. Find the following:

a) The prob. That a student scores over 80.

$$Y = \text{exam score}, Y \sim N(60, \sigma_Y = 10)$$

$$\Pr\{Y > 80\} = \Pr\left\{\frac{Y - \mu}{\sigma} > \frac{80 - 60}{10}\right\} = \Pr\{Z > 2\} = 1 - \Pr\{Z < 2\}$$
$$= 1 - 0.9772 = 0.0228$$

b) The prob. of scoring within two std dev's of the mean.

$$\Pr\{60 - 2(10) < Y < 60 + 2(10)\} = \Pr\{40 < Y < 80\} = \Pr\{Y < 80\} - \Pr\{Y < 40\}$$
$$= \Pr\{Z < \frac{80 - 60}{10}\} - \Pr\{Z < \frac{40 - 60}{10}\} = \Pr\{Z < 2\} - \Pr\{Z < -2\}$$
$$= 0.9772 - 0.0228 = 0.9544$$

c) The 90th percentile of exam score.

$$Z^{(90)} = 1.28 \quad (.90 \text{ is in row 1.2, column 0.08})$$

$$\text{So, } Y^{(90)} = \mu_Y + Z^{(90)}(\sigma_Y) = 60 + 1.28(10) = 72.8$$

I.e., 10% of students scored over 72.8, 90% below.