

AP Statistics Notes – Unit Six: Random Variable Distributions

Syllabus Objectives: 3.5 – The student will create probability distributions for discrete random variables, including geometric and binomial. 3.6 – The student will analyze probability distributions for discrete random variables, including geometric and binomial. 3.7 – The student will create probability distributions by using simulation techniques.

In the last unit, we learned that a “random phenomenon” was one that was unpredictable in the short term, but displayed a predictable pattern in the long run. In Statistics, we are often interested in numerical outcomes of random phenomena. In this unit, we will learn to define random variables to describe numerical outcomes of random phenomena as well as how to calculate the means and variances of such random variables.

- **Random Variables**

- **Random Variable** – A random variable is a variable whose value is a numerical outcome of a random phenomenon.
- **Discrete Random Variable** – A discrete random variable has a countable number of possible values.

Possible values of a
discrete random variable



- **Probability Distribution** – When describing a random variable “X”, be sure to note the probability distribution, showing the values X takes on and their respective probabilities. This can be done with a table or can be displayed using a **probability histogram**. The height of each bar represents the probability of the outcomes. The probabilities, p_i , must satisfy the following two conditions:
 1. $0 \leq p_i \leq 1$ for each i
 2. $p_1 + p_2 + \dots + p_k = 1$
- **Continuous Random Variable** – A continuous random variable takes on all possible values in an interval of numbers. We can display the probability distribution of a continuous random variable with a density curve. All continuous probability distributions assign a probability of zero to each individual outcome. Probabilities are defined over ranges of values.

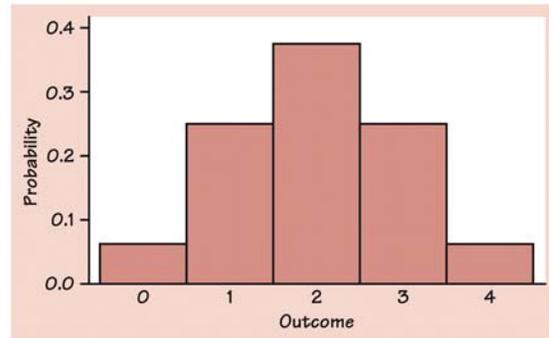
Possible values of a
continuous random variable



- **Random Variable Examples**

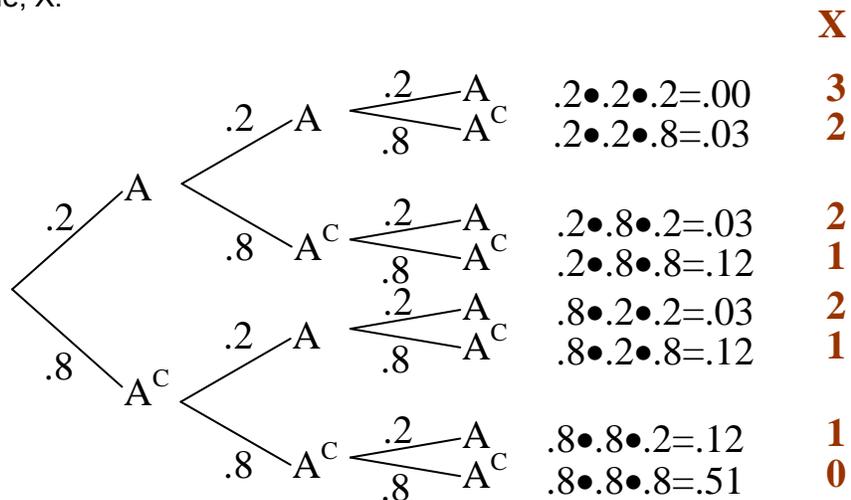
- **Descriptions of random variables**
 1. Experiment: A fair die is rolled. Random Variable: The number on the up face. Type: Discrete.
 2. Experiment: A pair of fair dice are rolled. Random Variable: The sum of the up faces. Type: Discrete.

- **Probability Histogram:**
The distribution can also be displayed as a histogram.



- **Question 1** – Find the probability of tossing at least two heads.
Solution: $P(X \geq 2) = 0.375 + 0.25 + 0.0625 = 0.6875$
 - **Question 2** – Find the probability of at least one head.
Solution: $P(X \geq 1) = 1 - P(X = 0) = 1 - 0.0625 = 0.9375$
- **Discrete Example 2** – Suppose that 20% of the apples sent to a sorting line are Grade A. If 3 of the apples sent to this plant are chosen randomly, determine the probability distribution of the number of Grade A apples in a sample of 3 apples.

- Consider the tree diagram to find the probabilities of each outcome of random variable, X.



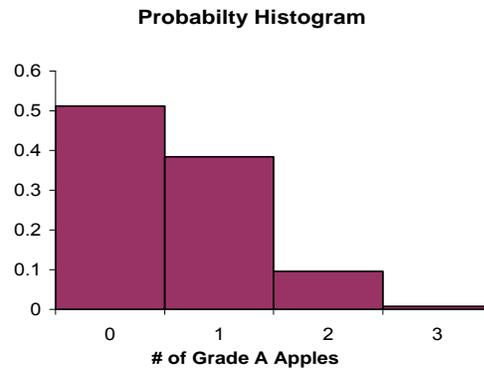
- The results in table form:

x	p(x)
0	$1(.8)^3$
1	$3(.8)^2(.2)^1$
2	$3(.8)^1(.2)^2$
3	$1(.2)^3$

or

x	p(x)
0	0.512
1	0.384
2	0.096
3	0.008

- The results in histogram form:

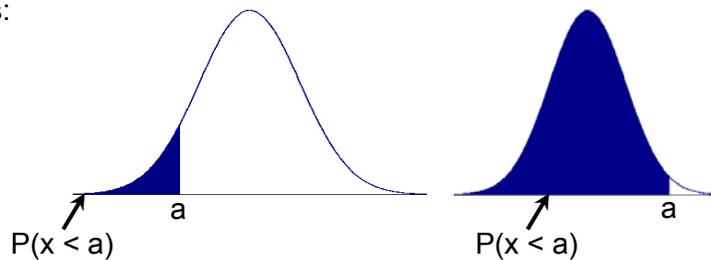


- **Continuous Random Variables**

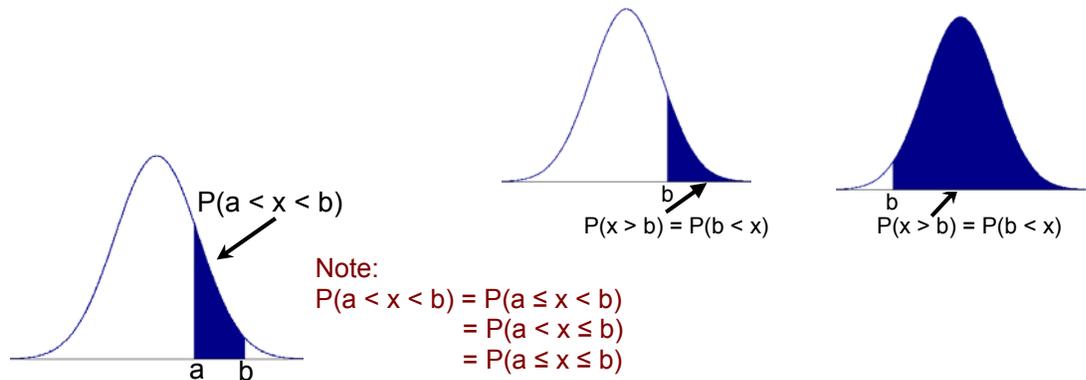
- A probability distribution for a continuous random variable, X , is specified by a mathematical function denoted by $f(x)$ which is called the density function. The graph of a density function is a smooth curve (the density curve). The following requirements must be met:
 1. $f(x) \geq 0$
 2. The total area under the density curve is equal to 1.

The probability that X falls in any particular interval is the area under the density curve that lies above the interval. We will use the density curves most familiar to us – rectangles, triangles and the normal curve.

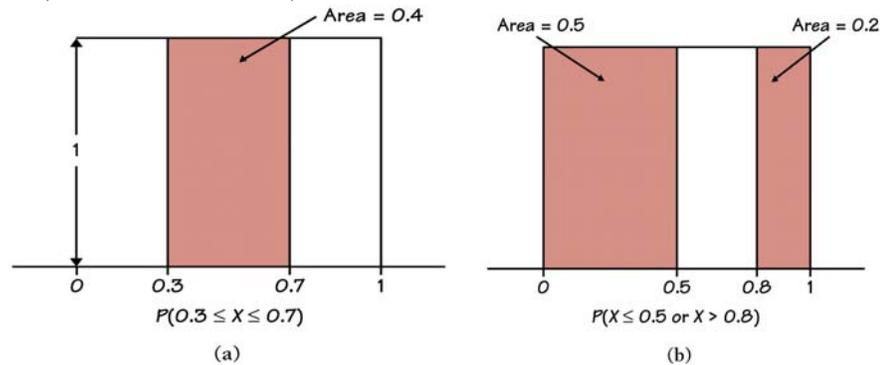
- Some illustrations:



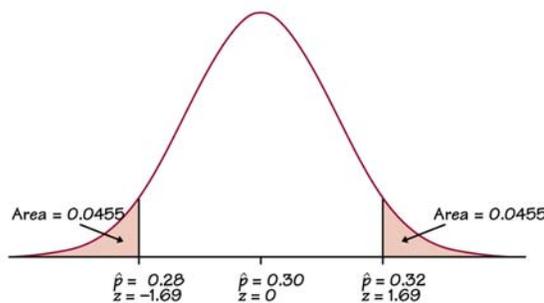
Notice that for a continuous variable, X , $P(x = a) = 0$ for any specific value a because the “area above a point” under the curve is a line segment and hence has no area. Specifically, this means $P(x < a) = P(x \leq a)$.



- **Continuous Example 1** – The random number generator will spread its output uniformly across the entire interval from 0 to 1 as we allow it to generate a long sequence of numbers. The results of many trials are represented by the density curve of a uniform distribution. The density curve has height 1 over the interval from 0 to 1.
 - **Question a** – Find the probability that the random number generator produces a number X between 0.3 and 0.7. $P(0.3 \leq X \leq 0.7)$
 - **Question b** - Find the probability that the random number generator produces a number X less than or equal to 0.5 and greater than 0.8. $P(X \leq 0.5 \text{ or } X > 0.8)$
 - **Solutions:** The height of the density curve is 1 and the area of a rectangle is the product of height and length, so the probability of any interval of outcomes is just the length of the interval. (a) $P(0.3 \leq X \leq 0.7) = 0.4$ and (b) $P(X \leq 0.5 \text{ or } X > 0.8) = 0.5 + 0.2 = 0.7$.



- **Continuous Example 2** – An opinion poll asks an SRS of 1500 American adults what they consider to be the most serious problem facing our schools. Suppose that if we could ask all adults the question, 30% would say “drugs”. Let us assume that this variable has an approximately normal distribution with a standard deviation of 0.0118. The mean of this distribution is 0.30. In this poll of 1500 adults, what is the probability that the poll result differs from the truth about the population by more than two percentage points?
 - **Solution:** Since this is a normal distribution, we have $N(0.3, 0.0118)$. The figure below shows this probability as an area under a normal density curve.



$$P(\hat{p} < 0.28 \text{ or } \hat{p} > 0.32) = P(\hat{p} < 0.28) + P(\hat{p} > 0.32)$$

$$P(\hat{p} < 0.28) = P\left(z < \frac{0.28 - 0.3}{0.0118}\right) = P(z < -1.69) = 0.0455$$

$$P(\hat{p} > 0.32) = P\left(z > \frac{0.32 - 0.3}{0.0118}\right) = P(z > 1.69) = 0.0455$$

$$P(\hat{p} < 0.28 \text{ or } \hat{p} > 0.32) = 0.0455 + 0.0455 = 0.0910$$

- **Binomial and Geometric Random Variables**

- A **binomial random variable** is a situation where these four conditions are satisfied:
 1. Each observation falls into one of just two categories, which for convenience we call “success” or “failure”.
 2. There is a fixed number n of observations.
 3. The n observations are all **independent**. That is, knowing the result of one observation tells you nothing about the other observations.
 4. The probability of success, call it p , is the same for each observation.

The binomial random variable, X , is defined as X = number of successes observed when experiment is performed. The probability distribution of X is called the binomial probability distribution.

- **The Binomial Formula:** If X has the binomial distribution with parameters n and p , the possible values of X are the whole numbers, $0, 1, 2, \dots, n$. The **binomial**

probability that X takes any value is: $P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$ The **binomial**

coefficient $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ counts the number of ways k successes can be arranged among n observations.

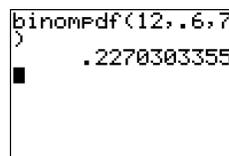
- **Binomial variable example 1** - The adult population of a large urban area is 60% black. If a jury of 12 is randomly selected from the adults in this area, what is the probability that precisely 7 jurors are black?

- **Solution:** Clearly, $n = 12$ and $p = .6$, so

$$P(X = 7) = \frac{12!}{7!5!} (.6)^7 (.4)^5 = 792(.02799)(.1024) = 0.2270. \text{ This can also be}$$

done on the TI-84, using the binomial probability distribution command. Use the following keystrokes:  and this will get you to the distribution

menu. Find the binompdf command and it is waiting for three parameters – n , p and X .



- Using the example above, if a jury of 12 is randomly selected from the adults in this area, what is the probability that less than 3 are black? This can also be solved using the binomial formula or the TI-84 using the binomcdf command.

▪ **Solution:**

$$\begin{aligned}
 P(x < 3) &= P(x < 2) = p(0) + p(1) + p(2) \\
 &= \frac{12!}{0!12!} (.6)^0 (.4)^{12} + \frac{12!}{1!11!} (.6)^1 (.4)^{11} + \frac{12!}{2!10!} (.6)^2 (.4)^{10} \\
 &= 0.00002 + 0.00031 + 0.00249 = 0.00281
 \end{aligned}$$

```
binomcdf(12,0.6,
2)
.0028101837
```

- **Binomial variable example 2** - On the average, 1 out of 19 people will respond favorably to a certain telephone sales pitch. If 25 people are called, what is the probability that at least two will respond favorably to this sales pitch?

▪ **Solution:**

$$P(X \geq 2) = 1 - P(X < 2) = 1 - P(X = 0) - P(X = 1)$$

$$1 - \frac{25!}{0!25!} \left(\frac{1}{19}\right)^0 \left(\frac{18}{19}\right)^{25} - \frac{25!}{1!24!} \left(\frac{1}{19}\right)^1 \left(\frac{18}{19}\right)^{24}$$

$$1 - 0.2588 - 0.3595 = 1 - 0.6183 = 0.3817$$

```
1-binomcdf(25,1/
19,1)
.3817436948
```

- A **geometric random variable** is a situation where these four conditions are satisfied:
 1. Each observation falls into one of just two categories, which for convenience we call “success” or “failure”.
 2. The variable of interest is the number of trials required to obtain the first success.
 3. The n observations are all **independent**. That is, knowing the result of one observation tells you nothing about the other observations.
 4. The probability of success, call it p , is the same for each observation.

The geometric random variable, X , is defined as X = number of trials until the first success is observed (including the success trial). The probability distribution of X is called the geometric probability distribution.

- **Computing the geometric variable:** If X has a geometric distribution with probability p of success and $(1 - p)$ of failure on each observation, the possible values of X are 1, 2, 3, ... If n is any one of these values, the probability that the first success occurs on the n th trial is: $P(X = n) = (1 - p)^{n-1} p$.

- $P(X > n)$ - The probability that it takes *more* than n trials to see the first success is:
 $P(X > n) = (1 - p)^n$.

- **Geometric variable example 1** - Over a very long period of time, it has been noted that on Friday's, 25% of the customers at the drive-in window at the bank make deposits. What is the probability that it takes 4 customers at the drive-in window before the first one makes a deposit?

- **Solution:** This problem is a geometric distribution as we are waiting for a success, where $p = 0.25$. To find this probability, we can use the formula or use the `geometpdf` command in the TI-84 calculator. Using the formula,
 $P(X = 4) = (0.75)^{4-1}(0.25) = (.421875)(0.25) = 0.1055$

To use the TI-84, go to the DISTR menu and find the `geometpdf` command. You need to input two parameters – p and X .

```
geometpdf(.25,4)
)
.10546875
```

- **Geometric variable example 2** - Roll a die until a 3 is observed. Find the following probabilities: a) The probability that it takes 4 rolls to observe a 3, b) The probability that it takes more than 6 rolls to observe a 3, c) The probability that we see a 3 before the 3rd roll.

- **Solution:** First, let's verify that this $X =$ the number of trials until a 3 occurs, is a geometric distribution. Note that rolling a 3 is a success and rolling any other number is a failure. The probability of rolling a 3 on each roll is the same – $1/6$. The observations are independent and we are rolling the die until a 3 appears. All of the 4 requirements are satisfied, so we can solve this

situation using a geometric distribution: a) $P(X = 4) = \left(\frac{5}{6}\right)^3 \left(\frac{1}{6}\right) = 0.0964$

b) $P(X > 6) = (1 - p)^n = \left(1 - \frac{1}{6}\right)^6 = \left(\frac{5}{6}\right)^6 = 0.3349$

c) $P(X \leq 2) = P(X = 1) + P(X = 2)$

$$= \left(\frac{5}{6}\right)^0 \left(\frac{1}{6}\right)^1 + \left(\frac{5}{6}\right)^1 \left(\frac{1}{6}\right)^1 = 0.1667 + 0.1389 = 0.3056 \text{ or}$$

$$P(X \leq 2) = 1 - P(X > 2) = 1 - \left(\frac{5}{6}\right)^2 = 1 - 0.6944 = 0.3056$$

```
geometpdf(1/6,4)
)
.0964506173
```

```
geometcdf(1/6,2)
)
.3055555556
```

```
1-geometcdf(1/6,
6)
)
.3348979767
```

Syllabus Objectives: 3.8 – The student will calculate the mean (expected value) and standard deviation of a random variable. 3.9 – The student will compute the mean and standard deviation for a linear transformation of a random variable.

- **Mean and standard deviation of a random variable**

- The **mean value of a discrete random variable, X**, denoted by μ_x , describes where the probability distribution of X is centered. This value is computed by first multiplying each possible X value by the probability of observing that value and then adding the resulting quantities.

$$\mu_x = \sum_{\text{all possible values of } x} x \cdot p(x)$$

- The **variance of a discrete random variable, X**, denoted by σ_x^2 , describes variability in the probability distribution. This value is computed by first subtracting the mean from each possible X value to obtain the deviations, then squaring each deviation and multiplying the result by the probability of the corresponding X value, and then finally adding these quantities.

$$\sigma_x^2 = \sum_{\text{all possible values of } x} (x - \mu_x)^2 \cdot p(x)$$

- The **standard deviation of a random variable, X**, denoted by σ_x , also describes variability in the probability distribution. The standard deviation is the square root of the variance.

$$\sigma_x = \sqrt{\sigma_x^2}$$

- **Example:** A professor regularly gives multiple choice quizzes with 5 questions. Over time, he has found the distribution of the number of wrong answers on his quizzes is as follows:

x	P(x)
0	0.25
1	0.35
2	0.20
3	0.15
4	0.04
5	0.01

- **Solution:** Multiple each X value by its probability and add the results to get

$$\mu_x =$$

x	P(x)	x•P(x)
0	0.25	0.00
1	0.35	0.35
2	0.20	0.40
3	0.15	0.45
4	0.04	0.16
5	0.01	0.05
		<u>1.41</u>

$$\mu_x = 1.41$$

x	P(x)	x•P(x)	x - μ	$(x - \mu)^2$	$(x - \mu)^2 \cdot P(x)$
0	0.25	0.00	-1.41	1.9881	0.4970
1	0.35	0.35	-0.41	0.1681	0.0588
2	0.20	0.40	0.59	0.3481	0.0696
3	0.15	0.45	1.59	2.5281	0.3792
4	0.04	0.16	2.59	6.7081	0.2683
5	0.01	0.05	3.59	12.8881	0.1289
		<u>1.41</u>			<u>1.4019</u>

$$\sigma_x^2 = 1.4019$$

$$\sigma_x = \sqrt{1.4019} = 1.184$$

- **Mean and standard deviation of a binomial random variable**

- The mean value (expected value) and the standard deviation of a binomial random variable are, respectively, $\mu_x = np$ and $\sigma_x = \sqrt{np(1-p)}$
- **Example:** A professor routinely gives quizzes containing 50 multiple choice questions with 4 possible answers, only one being correct. Occasionally, he just hands the students an answer sheet without giving them the questions and asks them to guess the correct answers. Let X be a random variable defined by X = number of correct answers on such an exam. Find the mean and standard deviation for X.

- **Solution:** The random variable satisfies all four requirements of a binomial distribution – there are only two outcomes, each question is independent, there is a set number of observations ($n = 50$) and the probability is constant ($p = 1/4$). This is binomial: $B(n, p) = B(50, 1/4)$. The mean and standard

deviation are: $\mu_x = np = 50\left(\frac{1}{4}\right) = 12.5$ and

$$\sigma_x = \sqrt{50\left(\frac{1}{4}\right)\left(\frac{3}{4}\right)} = \sqrt{9.375} = 3.06$$

- **Mean and standard deviation of a geometric random variable**

- If X is a geometric random variable with probability of success p on each trial, then the **mean**, or **expected value**, of the random variable, that is the expected number of trials required to get the first success, is $\mu_x = \frac{1}{p}$. The variance of X is

$$\sigma_x^2 = \frac{(1-p)}{p^2} \text{ with a standard deviation of } \sigma_x = \sqrt{\frac{(1-p)}{p^2}}.$$

- **Example:** A basketball player makes 80% of her free throws. We put her on the free-throw line and ask her to shoot free throws until she misses one. Let X = the number of free throws the player takes until she misses.
- **Solution:** This is a geometric setting as it satisfies all four conditions. The probability of this geometric variable is $p=0.80$. What is the expected number of free throws before she misses? We need to find the mean:

$\mu_x = \frac{1}{p} = \frac{1}{0.80} = 1.25$. To find the standard deviation of this distribution:

$$\sigma_x = \sqrt{\frac{(1-0.80)}{0.80^2}} = \sqrt{\frac{.20}{0.64}} = \sqrt{0.3125} = 0.559$$

- **Mean and standard deviation for a linear transformation of a random variable**

- If X is a random variable with mean μ_x and variance σ_x^2 and a and b are numerical constants, the random variable y defined by $y = a + bx$ is called a **linear function of the random variable X** .

- The mean of $y = a + bx$ is $\mu_y = \mu_{a+bx} = a + b\mu_x$.

- The variance of y is $\sigma_y^2 = \sigma_{a+bx}^2 = b^2\sigma_x^2$ and the standard deviation is

$$\sigma_y = \sigma_{a+bx} = b\sigma_x$$

- **Example:** Suppose X is the number of sales staff needed on a given day. If the cost of doing business on a day involves fixed costs of \$255 and cost per sales person per day is \$110, find the mean cost of doing business on a given day where the distribution of X is given below.

x	p(x)
1	0.3
2	0.4
3	0.2
4	0.1

- **Solution:** We need to find the mean of $y = 255 + 110x$

x	p(x)	xp(x)
1	0.3	0.3
2	0.4	0.8
3	0.2	0.6
4	0.1	0.4
		2.1

$$\mu_x = 2.1$$

$$\begin{aligned} \mu_y &= \mu_{255+110x} = 255 + 110\mu_x \\ &= 255 + 110(2.1) = \$486 \end{aligned}$$

We also need to find the variance and standard deviation of $y = 255 + 110x$.

x	p(x)	$(x-\mu)^2p(x)$
1	0.3	0.3630
2	0.4	0.0040
3	0.2	0.1620
4	0.1	0.3610
		0.8900

$$\sigma_x^2 = 0.89$$

$$\sigma_x = \sqrt{0.89} = 0.9434$$

$$\sigma_{255+110\mu_x}^2 = (110)^2 \sigma_x^2 = (110)^2 (0.89) = 10,769$$

$$\sigma_{255+110\mu_x} = (110)\sigma_x = (110)(0.9434) = 103.77$$

Syllabus Objectives: 3.10 – The student will determine the independence or dependence of two random variables. 3.11 – The student will calculate the mean and standard deviation for sums and differences of independent random variables.

- **Means and Variances for Combinations of random variables**

- Two random variables X and Y are **independent** if knowing that any event involving X alone did or did not occur tells us nothing about the occurrence of any event involving Y alone. Probability models often assume independence when the random variables describe outcomes that appear unrelated to each other. When random variables are not independent, the variance of their sum depends on the correlation between them and we will not be studying that case in this course. We only have to deal with the case of combining independent random variables.
- If X_1, X_2, \dots, X_n are random variables with means $\mu_1, \mu_2, \dots, \mu_n$ and variances $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$, and the variables are independent, the following formulas are used to find the mean and variances of their sums and differences.
 1. **Mean of a sum:** $\mu_{X_1+X_2+\dots+X_n} = \mu_1 + \mu_2 + \dots + \mu_n$
 2. **Mean of a difference:** $\mu_{X_1-X_2-\dots-X_n} = \mu_1 - \mu_2 - \dots - \mu_n$
 3. **Variance/S.D. of a sum:** $\sigma_Y^2 = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_n^2$ and $\sigma_Y = \sqrt{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_n^2}$
 4. **Variance/S.D. of a difference:** $\sigma_Y^2 = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_n^2$ and $\sigma_Y = \sqrt{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_n^2}$
- **Example 1:** A distributor of fruit baskets is going to put 4 apples, 6 oranges and 2 bunches of grapes in his small gift basket. The weights, in ounces, of these items are the random variables X_1, X_2 and X_3 respectively with means and standard deviations as given in the following table.

	Apple	Orange	Grape
Mean μ	8	10	7
Standard deviation σ	0.9	1.1	2

Find the mean, variance and standard deviation of the random variable Y= weight of fruit in a small gift basket.

- **Solution:** It is reasonable in this case to assume that the weights of the different types of fruit are **independent**.

$$\mu_Y = 4(\mu_1) + 6(\mu_2) + 2(\mu_3) = 4(8) + 6(10) + 2(7) = 106$$

$$\sigma_Y^2 = 4^2 (.9)^2 + 6^2 (1.1)^2 + 2^2 (2)^2 = 72.52 \text{ and } \sigma_Y = \sqrt{72.52} = 8.5159$$

- **Example 2:** Suppose “1 lb” boxes of Sugar Treats cereal have a weight distribution with a mean of $\mu_T = 1.050$ lbs and standard deviation $\sigma_T = 0.051$ lbs and “1 lb” boxes of Sour Balls cereal have a weight distribution with a mean $\mu_B = 1.090$ lbs and standard deviation $\sigma_B = 0.087$ lbs. If a promotion is held where the customer is sold a shrink wrapped package containing “1 lb” boxes of both Sugar Treats and Sour Balls cereals, what is the mean and standard deviation for the distribution of promotional packages?

- **Solution:** Combining these two values we get:

$$\mu_{T+B} = \mu_T + \mu_B = 1.05 + 1.09 = 2.14 \text{ lbs}$$

$$\sigma_{T+B}^2 = \sigma_T^2 + \sigma_B^2 = (0.051)^2 + (0.087)^2 = 0.002601 + 0.007569 = 0.01017$$

$$\sigma_{T+B} = \sqrt{0.01017} = 0.1008 \text{ lbs}$$

- **Example 3:** Tom and George are playing in the club golf tournament. Their scores vary as they play the course repeatedly. Tom’s score X has the $N(110,10)$ distribution, and George’s score Y varies from round to round according to the $N(100,8)$ distribution. If they play independently, what is the probability that Tom will score lower than George and then do better in the tournament?

- **Solution:** The difference $X - Y$ between their scores is normally distributed. We need to find the mean and variance of this distribution.

$$\mu_{X-Y} = \mu_X - \mu_Y = 110 - 100 = 10$$

$$\sigma_{X-Y}^2 = \sigma_X^2 + \sigma_Y^2 + 10^2 + 8^2 = 164$$

Because $\sqrt{164} = 12.8$, $X - Y$ has the $N(10, 12.8)$ distribution.

$$P(X < Y) = P(X - Y < 0)$$

$$= P\left(\frac{(X - Y) - 10}{12.8} < \frac{0 - 10}{12.8}\right)$$

$$= P(Z < -0.78) = 0.2177$$

