

STT 315

Preparation for the FINAL EXAM (longer than any actual exam).

Note: According to the syllabus your course grade is based upon total points. Each graded activity such as exams, homework, bonus, is converted into such points by means of a curve. The relationship of total course points to course grade is given in the syllabus.

New: Your total course points will be **the larger of**

- Total points figured by the method outlined in the syllabus (usual method).
- Total points figured as in (a), except that the points you earn for your lowest exam will be replaced by **(23/30) times (points you earn for the final exam)**.

1. $P(A) = 0.7, P(B) = 0.4$.

a. What is the very least $P(AB)$ could be?

ans. a. $P(A \cup B)$ is the same as $P(A \text{ or } B)$ and cannot exceed one. But

b. $P(A \text{ or } B) = P(A) + P(B) - P(AB) = 0.7 + 0.4 - P(AB) = 1.1 - P(AB)$
so $P(AB)$ must be at least 0.1. If $P(AB) = 0.1$ then $P(A \text{ or } B) = 1.1 - 0.1 = 1$. For this case

c. $P(B | A) = P(AB) / P(A) = 0.1 / 0.7$

d. $P(AB) = 0.1$ as shown above. AB is another name for 'A and B.'

e. $P(A^c) = 1 - P(A) = 1 - 0.7 = 0.3$

b. If A, B are independent events what are

a. $P(AB) = P(A) P(B) = 0.7 \cdot 0.4 = 0.28$

b. $P(A \text{ or } B) = P(A) + P(B) - P(AB) = 0.7 + 0.4 - 0.28$

c. $P(B | A) = P(B) = 0.4$ since, for independent events, knowing A has occurred does not change the probability for B. Formally, $P(AB) / P(A) = ($ for independent events)
 $P(A) P(B) / P(A) = P(B) = 0.4$.

d. $P(B | A^c) = P(B) = 0.4$.

2. Let events A, B be

A = customer will buy the product

B = customer leaves a cash deposit

Suppose that $P(A) = 0.3, P(B | A) = 0.8, P(B | A^c) = 0.4$.

a. Make a complete tree diagram for this information.

$$\begin{array}{l} P(A) = 0.3 \\ \quad P(B|A) = 0.8 \quad P(AB) = P(A) P(B|A) = 0.3 \cdot 0.8 = 0.24 \end{array}$$

$$\quad P(B^c|A) = 0.2 \quad P(AB^c) = P(A) P(B^c|A) = 0.3 \cdot 0.2 = 0.06$$

$$\begin{array}{l} P(A^c) = 0.7 \\ \quad P(B|A^c) = 0.4 \quad P(A^cB) = P(A^c) P(B|A^c) = 0.7 \cdot 0.4 = 0.28 \end{array}$$

$$\quad P(B^c|A^c) = 0.6 \quad P(A^cB^c) = P(A^c) P(B^c|A^c) = 0.7 \cdot 0.6 = 0.42$$

b. Make a complete Venn diagram for this information.



c. Determine $P(B)$ and also $P(B^c)$.

$$P(B) = P(AB) + P(A^c B) = 0.24 + 0.28 = 0.52 \text{ (corrected)}$$

$$P(B^c) = 1 - P(B) = 1 - 0.52 = 0.48 \text{ (corrected)}$$

d. Determine $P(A | B)$ and also $P(A | B^c)$.

$$P(A|B) = P(AB) / P(B) = 0.24 / 0.52 \text{ (corrected)}$$

3. $P(OIL) = 0.2$, $P(+ | OIL) = 0.8$, $P(+ | OIL^c) = 0.3$, cost to test is 40, cost to drill is 300, return from oil is 1200.

a. Make a complete tree diagram including all endpoint probabilities and the consequent net returns x from the policy “test, but drill only if the test is positive.”

	$P(+ OIL) = 0.8$	$P(OIL+) = P(OIL) P(+ OIL) = 0.2 \cdot 0.8 = 0.16$	net 860
$P(OIL) = 0.2$	$P(- OIL) = 0.2$	$P(OIL-) = P(OIL) P(- OIL) = 0.2 \cdot 0.2 = 0.04$	-40

	$P(+ OIL^c) = 0.3$	$P(OIL^c+) = P(OIL^c) P(+ OIL^c) = 0.8 \cdot 0.3 = 0.24$	-340
$P(OIL^c) = 0.8$	$P(- OIL^c) = 0.7$	$P(OIL^c-) = P(OIL^c) P(- OIL^c) = 0.8 \cdot 0.7 = 0.56$	-40

b. Find $E X$.

The net returns x , and their probabilities, for each of the above for tree endpoints are

	x	$p(x)$	$x p(x)$
$-40 - 300 + 1200 =$	860	0.16	860 0.16
$-40 - 00 + 00 =$	-40	0.04	-40 0.04
$-40 - 300 + 00 =$	-340	0.24	-340 0.24
$-40 - 00 + 00 =$	-40	0.56	-40 0.56

$$E X = \text{sum of } x p(x) = 32$$

c. Find $E Y$ where y is the net return from the policy “just drill, do not test.”

$$E Y = 0.2 (1200 - 300) + 0.8 (-300) = -60.$$

4. Drawing balls with equal probability but without replacement from $\{R R G G G Y\}$.

a. $P(Y2)$ (guess it from a principle that you properly name and confirm your guess using the rules of total probability and multiplication). Identify your use of the rules.

$$\begin{aligned}
 \text{Order of the deal does not matter so } P(Y2) &= P(Y1) = 1/6. \text{ Using total probability we} \\
 \text{have } P(Y2) &= P(Y1 Y2) + P(Y1^c Y2) = 0 + P(Y1^c Y2) = P(Y1^c) P(Y2 | Y1^c) \\
 &= 5/6 \cdot 1/5 = 1/6
 \end{aligned}$$

b. $P(G1 | Y2)$. Prove your answer using the rules.

$P(G1 | Y2)$ is the same as $P(G2 | Y1) = 3 / 5$ since **order of the deal does not matter**. To prove it using the rules use $P(G1 | Y2) = P(G1 Y2) / P(Y2)$ where

$$P(G1 Y2) = P(G1) P(Y2 | G1) = 3/6 \cdot 1/5$$

$$P(Y2) = P(Y1) = 1/6$$

So $P(G1 | Y2) = (3/6 \cdot 1/5) / (1/6) = 3/5$, confirming the fact that order of the deal does not matter.

5. $P(\text{rain today}) = 0.6$, $P(\text{rain tomorrow}) = 0.5$, $P(\text{rain tomorrow} | \text{rain today}) = 0.8$.

a. $P(\text{rain today and rain tomorrow})$.

$$P(\text{tod tom}) = P(\text{tod}) P(\text{tom} | \text{tod}) = 0.6 \cdot 0.8 = 0.48$$

b. Give the complete Venn diagram.

$$\text{tod tom} \quad 0.6 \cdot 0.8 = 0.48$$

$$\text{tod tom}^c \quad 0.6 (1 - 0.8) = 0.12$$

using the decomposition $\text{tod}^c \text{tom} = \text{tom} - \text{tod tom}$ (draw a picture of this!)

$$\text{tod}^c \text{tom} \quad P(\text{tom}) - P(\text{tod tom}) = 0.5 - 0.6 \cdot 0.8 = 0.02$$

$$\text{tod}^c \text{tom}^c \quad 1 - \text{sum of above} = 1 - (0.48 + 0.12 + 0.02) = \mathbf{0.38}$$

c. Give the complete tree diagram.

$$P(\text{tom} | \text{tod}) = 0.8$$

$$P(\text{tod tom}) = 0.6 \cdot 0.8 = 0.48$$

$$P(\text{tod}) = 0.6$$

$$P(\text{tom}^c | \text{tod}) = 0.2$$

$$P(\text{tod tom}^c) = 0.6 \cdot 0.2 = 0.12$$

$$P(\text{tom} | \text{tod}^c) = \mathbf{0.02} / 0.4$$

$$P(\text{tod}^c \text{tom}) = 0.5 - 0.6 \cdot 0.8 = \mathbf{0.02}$$

$$P(\text{tod}^c) = 0.4$$

$$P(\text{tom}^c | \text{tod}^c) = \mathbf{0.38} / 0.4$$

$$P(\text{tod}^c \text{tom}^c) = 1 - 0.62 = \mathbf{0.38}$$

d. Give $P(\text{rain today} | \text{rain tomorrow})$.

$$P(\text{tod} | \text{tom}) = P(\text{tod tom}) / P(\text{tom}) = 0.48 / 0.5$$

6.	x	p(x)	x p(x)	$x^2 p(x)$
	0	0.2	0 · 0.2 = 0.0	0 · 0.2 = 0.0
	1	0.8	1 · 0.8 = 0.8	1 · 0.8 = 0.8
sums			$E X = 0.8$	$E X^2 = 0.8$

a. Give the formula for $E X$ for a r.v. taking values 0 or 1 only, in terms of $p = P(X = 1)$.

$E X = p$ for a r.v. taking only values 0, 1 with respective probabilities p, q = 1 - p.

So $E X = 0.8$ (confirmed above).

b. Calculate $E X = \text{“sum of } x p(x)\text{”}$ confirming (a).

Done above.

c. Give the formulas for Var X and sd X (for a r.v. taking values 0 or 1 only) in terms of $p = P(X = 1)$ and $q = P(X = 0)$.

$$\text{Var } X = pq = 0.8 \cdot 0.2 = 0.16. \quad \text{Sd } X = \text{root}(pq) = \text{root}(0.16) = 0.4.$$

d. Calculate Var X and sd X from values of $\{x, p(x)\}$ confirming your answer (c).

$E X = 0.8 = E X^2$ as shown in the table above. So $\text{Var } X = E X^2 - (E X)^2 = 0.8 - 0.64 = 0.16$. So $\text{sd } X = \text{root}(0.16) = 0.4$. This confirms (c).

7.	x	p(x)	x p(x)	$x^2 p(x)$	$(x - E X)^2 p(x)$
	3	0.4	$3 \cdot 0.4 = 1.2$	$9 \cdot 0.4 = 3.6$	$(3 - 1.0)^2 \cdot 0.4 = 1.6$
	-1	0.2	$-1 \cdot 0.2 = -0.2$	$1 \cdot 0.2 = 0.2$	$(-1 - 1.0)^2 \cdot 0.2 = 0.8$
	0	0.4	$0 \cdot 0.4 = 0.0$	$0 \cdot 0.4 = 0.0$	$(0 - 1.0)^2 \cdot 0.4 = 0.4$
sums			$E X = 1.0$	$E X^2 = 3.8$	$\text{Var } X = 2.8$

a. Calculate $E X = 1.0$ above.

b. Calculate Var X and sd X using the **definition** (as opposed to the computing formula). $\text{Var } X = 2.8$ as in the last column above. $\text{Sd } X = \text{root}(\text{Var } X) = \text{root}(2.8)$.

c. Re-calculate Var X and sd X using the **computing formula**, confirming (b).

$\text{Var } X = \text{expected square of } X - \text{square of expected } X = 3.8 - 1^2 = 2.8$ from above, confirming (b).

d. R.v. $Y = 3 X - 2$. Determine $E Y$, $\text{Var } Y$ and $\text{sd } Y$ from your answers above, exploiting known connections with $E X$, $\text{Var } X$ and $\text{sd } X$.

$$E Y = 3 E X - 2 = 3(1.0) - 2 = 1.0 \text{ also.}$$

$$\text{Var } Y = \text{Var}(3 X - 2) = \text{Var}(3 X) = 9 \text{ Var } X = 9(2.8) = 25.2.$$

$$\text{Sd } Y = \text{root}(\text{Var } Y) = \text{root}(25.2) = 3 \text{ root}(2.8) = 3 \text{ sd}(X).$$

e. Re-calculate $E Y$, $\text{Var } Y$ and $\text{sd } Y$ directly from the table below confirming (d).

($p(x) = p(y)$ since each y corresponds to precisely one x)

	y	x	p(x)	y p(y)	$y^2 p(y)$
	7	3	0.4	$7 \cdot 0.4 = 2.8$	$49 \cdot 0.4 = 19.6$
	-5	-1	0.2	$-5 \cdot 0.2 = -1.0$	$25 \cdot 0.2 = 5.0$
	-2	0	0.4	$-2 \cdot 0.4 = -0.8$	$4 \cdot 0.4 = 1.6$
sums				$E Y = 1.0$	$E Y^2 = 26.2$

So $E Y = 1.0$ and $\text{Var } Y = E Y^2 - (E Y)^2 = 26.2 - 1.0^2 = 25.2$, both as in (d).

8. R.v. X and Y have

$$E X = 6$$

$$E Y = -3$$

$$\text{Var } X = 9$$

$$\text{Var } Y = 13$$

a. Determine $E(2 X + 3 Y - 4) = 2 E X + 3 E Y - 4 = 2(6) + 3(-3) - 4 = -1$.

b. If X, Y are independent determine $\text{Var}(2 X + 3 Y - 4)$ showing how independence is used.

$$\begin{aligned}\text{Var}(2 X + 3 Y - 4) &= \text{Var}(2 X + 3 Y) = (\text{if independent}) \text{Var}(2 X) + \text{Var}(3 Y) \\ &= 4 \text{Var } X + 9 \text{Var } Y = 4 \cdot 9 + 9 \cdot 13 = 153.\end{aligned}$$

c. From (b) determine $\text{sd}(2 X + 3 Y - 4) = \sqrt{153}$

9. CLT. Each account of a population of business accounts is scored with x = balance due. Suppose the population mean balance is due is $E X = \$466.48$ population sd $\sigma = \$74.88$. A with replacement sample of 100 accounts will be selected.

a. Are you able to sketch the population distribution of x ?

No. We know only the mean and sd of the population. There is always the possibility of a population “in control” in which case the distribution is normal with mean 466.48 and sd 74.88. However, there are infinitely many other distributions having these very same mean and sd.

b. Determine $E \bar{x}$ and $\text{sd } \bar{x}$.

$E \bar{x} = \text{population mean} = 466.48$. That is, the sample mean, averaged over all of its possibilities according to their probabilities, is exactly the mean of the population from which the samples are being selected. This is true for with-replacement sampling, as here, but is also true for without-replacement sampling.

$$\text{sd } \bar{x} = \sigma / \sqrt{n} = 74.88 / \sqrt{100} = 7.488.$$

c. Sketch the approximation of the distribution of \bar{x} offered by the central limit theorem (CLT) identifying $E \bar{x}$ and $\text{sd } \bar{x}$ as recognizable elements of your sketch.

Bell curve (normal density) having mean 466.48 and sd 7.488.

d. Repeat (c) if instead the sample is selected without replacement and the population size is $N = 4000$. Is your sketch much different from (c) in this case?

Sampling without-replacement alters $\text{sd } \bar{x}$ to

$$7.488 \text{ times FPC} = 7.488 \sqrt{\frac{4000-100}{4000-1}} = 7.3947 \text{ (little changed)}$$

10. CLT. A population of accounts has 30% that are overdue. A sample of 400 accounts is to be selected with-replacement.

a. Sketch the approximate distribution of \hat{p} (fraction of overdue accounts in the sample). Clearly identify $E \hat{p}$ and $sd \hat{p}$ as recognizable elements of your sketch and evaluate them numerically.

$$\hat{p} = (\text{number of overdue accounts in the sample of 400}) / 400$$

(it is random since the sample is random).

$E \hat{p}$ = population fraction of overdue accounts = $p = 0.3$ (we're told 30% are overdue).

$$sd \hat{p} = \sqrt{pq / 400} = \sqrt{(0.3 \cdot 0.7) / 400} \sim 0.023.$$

Ans. Bell curve with mean 0.3 and sd 0.023.

The ACTUAL margin of error for \hat{p} is therefore $\pm 1.96 \cdot 0.023 \sim 0.045$.

Had we not been told $p = 0.3$ we would estimate ME from the sample as

$$ME \sim \pm 1.96 \sqrt{\hat{p}\hat{q} / 400} \text{ (often referred to as margin of error)}$$

A 95% z-based CI for p would be $\hat{p} \pm 1.96 \sqrt{\hat{p}\hat{q} / 400}$.

b. Repeat (a) except assume that the sample is without-replacement and the population size is 3000. Is this sketch much different from (a)?

$$FPC = \sqrt{(3000 - 400) / (3000 - 1)} \sim 0.93.$$

Bell curve with mean $p = 0.3$ and $sd = 0.023 \cdot FPC = 0.023 \cdot 0.93$.

11. We average around 6.4 shortages in a week. The Poisson distribution is thought to apply to X = number of shortages in one week.

a. Determine the probability that there are 5 shortages in one week.

$$p(x) = e^{-\mu} \mu^x / x! = e^{-6.4} 6.4^5 / 5! \sim 0.149.$$

b. Out of 52 weeks, assuming this model applies to all weeks, around how many weeks should experience exactly 5 shortages?

$$52 \text{ times } 0.149 = 7.73 \sim 8.$$

c. Sketch the normal approximation of the distribution of X . Clearly identify $E X$ and $sd X$ (which you evaluate numerically) as recognizable elements of your sketch.

$$E X = 6.4 > 3. \text{ Sd of Poisson} = \sqrt{\text{mean}} = 2.53.$$

Bell curve having mean 6.4 and sd 2.53 (a fairly good approximation for Poisson with mean > 3).

d. Use the z-table and (c) to approximate $P(X \text{ in range } 3, 4, 5, 6)$ using the continuity correction. Identify the two z-scores you are using.

$$z = (2.5 - 6.4) / 2.53 = -1.54$$

(continuity correction dips back to 2.5 to capture more of $p(3)$ under z-curve)

$$z = (6.5 - 6.4) / 2.53 = 0.04$$

$$p(3) + p(4) + p(5) + p(6) \sim P(-1.54 < Z < 0.04) = P(0 < Z < 1.54) + P(0 < Z < 0.04)$$

$$0.4382 + 0.016 = 0.4398$$

The actual values are $p(3) + p(4) + p(5) + p(6)$

$$= 0.0725945 + 0.116151 + 0.148674 + 0.158585 = 0.4960.$$

e. Use closest z-entry to approximate the 75th percentile of X (you first need to find the 75th percentile of z then convert to an x-score).

75th percentile of Z captures 25% between 0 and itself. So we seek z with $P(0 < Z < z) = 0.25$. This is the reverse table use. We enter 0.25 to the body of the z-table finding $z = 0.67$ corresponds to the closest entry 0.2486.

12. A with replacement sample of 50 customers, for score $x =$ dollar value returns they have made to us last year, finds sample mean $\bar{x} = \$4.89$ with sample sd $s = \$2.10$.

a. Estimate the population mean and population sd.

estimate of pop mean is $\bar{x} = 4.89$.

estimate of pop sd is $s = 2.10$

b. Estimate the sd of \bar{x} .

theory says sd of \bar{x} is σ / \sqrt{n}

we estimate sd of \bar{x} by $s / \sqrt{n} = 2.10 / \sqrt{50} = 0.297$.

c. Determine a 95% z-based CI for $\mu =$ population mean.

usual $\bar{x} \pm 1.96 s / \sqrt{n} = 4.89 \pm 1.96 2.10 / \sqrt{50} = [4.30791, 5.47209]$.

$P(\mu \text{ in (random) CI } \bar{x} \pm 1.96 s / \sqrt{n}) \sim 0.95$.

13. A **normal** population is sampled finding $\{3.4, 3.8, 4.8\}$.

a. Calculate \bar{x} and s from this data.

$$\bar{x} = (3.4 + 3.8 + 4.8) / 3 = 4$$

$$s = \sqrt{((3.4 - 4)^2 + (3.8 - 4)^2 + (4.8 - 4)^2) / (3 - 1)} = \sqrt{0.52} = 0.721.$$

b. Estimate the population mean and sd from this data.

est of μ is $\bar{x} = 4$

est of σ is $s = 0.721$

c. Estimate sd \bar{x} from this data.

est of sd of \bar{x} is $s / \sqrt{n} = 0.721 / \sqrt{3} = 0.416$

d. Determine a 90% t-based CI for the population mean μ from this data.

$\bar{x} \pm t s / \sqrt{n} = 4 \pm 2.92 0.721 / \sqrt{3} = [2.78449, 5.21551]$.

DF = $3 - 1 = 2$.

e. What is the value of $P(\mu \text{ in } 90\% \text{ t-based CI}) = 0.9$
(**exact** if we use t and data calculated to infinite precision).

f. To what sample size n_{FINAL} must we continue in order to obtain a 90% t-based hybrid CI of the form $\bar{x}_{\text{FINAL}} \pm 0.2$?
formula $n \sim (t_{s/B})^2 = (2.92 \cdot 0.721 / 0.2)^2 \sim 111$.

Since this n is so large, it might be best to just toss out the original $n_0 = 3$ samples and sample 111 fresh, then use a z-based CI. We'll not pursue any detailed comparison of the two methods except to say it may be best whenever the indicated n_{FINAL} is large.

g. If we do continue sampling to the n_{FINAL} of part (f) finding $\bar{x}_{\text{FINAL}} = 4.22$ what will be our hybrid t-based CI from (f)?
 $\bar{x}_{\text{FINAL}} \pm 0.2 = 4.22 \pm 0.2$ (really, it is almost exactly this)

14. We desire a 95% z-based CI for $p =$ fraction of viewers of our advertising who have a "favorable impression of our company." A with-replacement sample of 100 viewers finds that 61 have a favorable impression.

a. Give the 95% z-based CI for p based on the above data.
 $p_{\text{HAT}} \pm 1.96 \text{ root}(p_{\text{HAT}} q_{\text{HAT}} / n) = 0.61 \pm 1.96 \text{ root}(0.61 \cdot 0.39 / 100)$
 $= [0.514401, 0.705599]$.

b. Determine an n_{FINAL} to which we must continue sampling in order to achieve a 95% z-based CI of the form $p_{\text{HAT}} \pm 0.1$.
 $n \sim (z \text{ root}(p_{\text{HAT}} q_{\text{HAT}}) / B)^2 = (1.96 \text{ root}(0.61 \cdot 0.39) / 0.1)^2 \sim 92$

c. If we do continue sampling to n_{FINAL} finding 63.7% of them have a favorable impression of our company give the resulting 95% hybrid z-based CI for p .

Just quote the CI (a) for the data you have, we have the desired precision, and even a little more, already. Had n_{FINAL} exceeded 100 we'd have used hybrid CI $p_{\text{HATfinal}} \pm 0.1$

15. It is desired to estimate the population average amount our typical corporate business client will spend with us next quarter. We sample 100 accounts with-replacement and (carefully) spend some time and money evaluating their likely purchase needs y next quarter. For each account we also have the score $x =$ amount they spent this quarter (on file). From the 100 we find

$$\begin{aligned} \bar{x}_{\text{BAR}} &= \$342000 & \bar{y}_{\text{BAR}} &= \$287000 \\ s(x) &= \$88000 & s(y) &= \$72000 \\ \text{sample correlation } \rho_{\text{HAT}} &= 0.85 \end{aligned}$$

Suppose it is known that $\mu_x = 338799$ (i.e. avg of all accounts this quarter)

a. Give a 95% z-based CI for $\mu(y)$ based on the y -data alone.
 $\bar{y}_{\text{BAR}} \pm 1.96 s_y / \text{root}(n) = 287000 \pm 1.96 72000 / \text{root}(100) = [272888, 301112]$.

b. Give the regression estimate of $\mu(y)$. Show that it differs from $y\text{BAR}$ used in (a).
 $y\text{BAR} + (\text{mux} - x\text{BAR}) (s_y / s_x) \text{rhoHAT}$
 $= 287000 + (338799 - 342000) (72000 / 88000) 0.85 = 284774$ (less than $y\text{BAR}$)

Rationale: Since $x\text{BAR}$ has over-estimated mux (known), and x is positively correlated with y , we reduce $y\text{BAR}$.

c. Give the estimate of sd (b) showing that it differs from the estimate of $\text{sd}(y\text{BAR})$.
 The estimated sd of the regression-based estimator (b) of $\mu(y)$ is
 $(s_y / \text{root}(n)) \text{root}(1 - \text{rhoHAT}^2) = (72000 / \text{root}(100)) \text{root}(1 - 0.85^2) = 3793$.
 This is less than the estimated sd of $y\text{BAR}$ which is 7200. It is as though instead of 100 samples we'd had $100 (7200 / 3793)^2 \sim 360$. That is a big advantage held by the regression method. However, you'd need to know mux and sample x -scores.

d. Give the 95% z -based regression-estimator based CI for $\mu(y)$ showing that it is narrower than (a) for the same sample size 100.
 $\text{regr est} \pm 1.96 (s_y / \text{root}(n)) \text{root}(1 - \text{rhoHAT}^2) = 284774 \pm 3793 = [280981, 288567]$.
 This CI is far narrower than the usual CI (a) based upon $y\text{BAR}$ alone. To use the regression based approach we have to know mux and be able to pick up x -scores along with y scores, but the savings can be considerable if rhoHAT is near -1 or 1 .

e. Modify (d) if the sample is without-replacement and the population size is 3000.
 $\text{regr est} \pm 1.96 (s_y / \text{root}(n)) \text{root}(1 - \text{rhoHAT}^2) \text{FPC}$
 $= 284774 \pm 1.96 3793 \text{root}((3000-100) / (3000-1)) = [281044, 288504]$
 It has $\text{FPC} = 0.983356$ times width of the regr -based CI for with-replacement sampling.

16. **Independent** samples of 60 women and 40 men, all professionals living alone, are drawn with-replacement finding

women	$x\text{BAR} = 24$	$s(x) = 15$	$n(x) = 60$
men	$y\text{BAR} = 18$	$s(y) = 10$	$n(y) = 40$

where x = amount spent dining out, y = amount spent dining out.

a. Give an estimate of $\mu(x) - \mu(y)$.

A usual estimator is $x\text{BAR} - y\text{BAR} = 24 - 18 = 6$.

In place of $x\text{BAR}$ we sometimes use more sophisticated estimators, not discussed here, such a “toss out the very largest and smallest of the data then average the rest.” Such alternative estimators are used to guard against “outliers” (e.g. Bill Gates throws off an average income if he is included but the estimator just described would leave him out).

b. Give estimates of (population) sigmaWOMEN and sigmaMEN .

est of sigmaWOMEN is $s\text{WOMEN} = 15$

est of sigmaMEN is $s\text{MEN} = 10$.

c. Give estimates of $\text{sd } x\text{BAR}$, $\text{sd } y\text{BAR}$, $\text{sd}(x\text{BAR} - y\text{BAR})$.

est of sd of $x\text{BAR}$ is $s(x) / \text{root}(n_x) = 15 / \text{root}(60) = 1.94$

est of sd of $y\text{BAR}$ is $s(y) / \text{root}(n_y) = 10 / \text{root}(40) = 1.58$

est of sd of $(x\text{BAR} - y\text{BAR})$ is $\text{root}(s(x)^2 / n(x) + s(y)^2 / n(y))$

$= \text{root}(225 / 60 + 100 / 40) = 2.5$ (more variable than each of $x\text{BAR}$, $y\text{BAR}$ alone)

d. Give 95% z-based CI for $\mu(x) - \mu(y)$.

$(\bar{x} - \bar{y}) \pm 1.96 \text{ est sd of } (\bar{x} - \bar{y}) = 6 \pm 1.96 \cdot 2.5 = [1.1, 10.9]$
which is uncomfortably wide compared with the estimate $\bar{x} - \bar{y} = 6$.

17. **Independent** samples of 60 women and 40 men, all working the same job, are drawn with-replacement finding

women 35 have hand cramping

men 17 have hand cramping

Define $p(x)$ = fraction of all women in this job who have hand cramping, $p(y)$ the corresponding rate for men.

a. Estimate each of $p(x)$, $p(y)$.

$$\hat{p}_x = 35 / 60 = 0.58333$$

$$\hat{p}_y = 17 / 40 = 0.425$$

b. Give estimates of $\text{sd } \hat{p}_x$, $\text{sd } \hat{p}_y$ and $\hat{p}_x - \hat{p}_y$.

$$\text{est of sd of } \hat{p}_x \text{ is } \sqrt{\hat{p}_x \hat{q}_x / n_x} = \sqrt{.58333 \cdot .41667 / 60} = 0.064$$

$$\text{est of sd of } \hat{p}_y \text{ is } \sqrt{\hat{p}_y \hat{q}_y / n_y} = \sqrt{.425 \cdot .575 / 40} = 0.078$$

$$\text{est of sd of } (\hat{p}_x - \hat{p}_y) \text{ is } \sqrt{\hat{p}_x \hat{q}_x / n_x + \hat{p}_y \hat{q}_y / n_y} \\ = \sqrt{.58333 \cdot .41667 / 60 + .425 \cdot .575 / 40} = 0.101 \text{ (larger than for either } \hat{p}_x)$$

c. Give a 90% z-based CI for $p(x) - p(y)$.

$$(\hat{p}_x - \hat{p}_y) \pm 1.645 \text{ est sd of } (\hat{p}_x - \hat{p}_y)$$

$$(.58333 - .425) \pm 1.645 \cdot 0.101 = [-0.007815, 0.324475].$$