

Statistics 23, Section 1, Final Exam
Tuesday, December 14, 1999

Name: _____

Pledge: I have neither given nor received aid on this examination.

Signature: _____

Instructions: Show all work, but do not do hard arithmetic (an answer like $8.3 + \frac{7}{\sqrt{5.1}}$ is fine).

1. An experiment results in one of three mutually exclusive outcomes, A , B or C . It is known that $P(A) = 0.2$, $P(B) = 0.4$ and $P(C) = 0.3$.

a. Find $P(A \text{ or } B)$

[5]

$$\begin{aligned} P\{A \text{ or } B\} &= P\{A\} + P\{B\} - P\{A \text{ and } B\} \\ &= 0.2 + 0.4 - 0 \quad (\text{since mut. exc.}) \\ &= 0.6 \end{aligned}$$

b. Find $P(B \text{ and } C)$

[5]

$$0 \quad (\text{since mut. exc.})$$

c. Find $P(A|C)$

[5]

$$0 \quad (\text{since mut. exc.})$$

d. Are A and B independent? Why or why not?

[5]

No, knowing A occurs means B does not occur, so have changed chances
i.e. $P\{B|A\} = 0$, not equal to $P\{B\}$

2. To study the effects of competition on cable television rates, 4 counties were selected at random, and their rates before and after the introduction of competition were put in an Excel spreadsheet as:

	A	B	C
1	County	Rate Before:	Rate After:
2	A	\$21.35	\$21.56
3	B	\$25.73	\$25.91
4	C	\$18.92	\$19.32
5	D	\$22.07	\$22.35

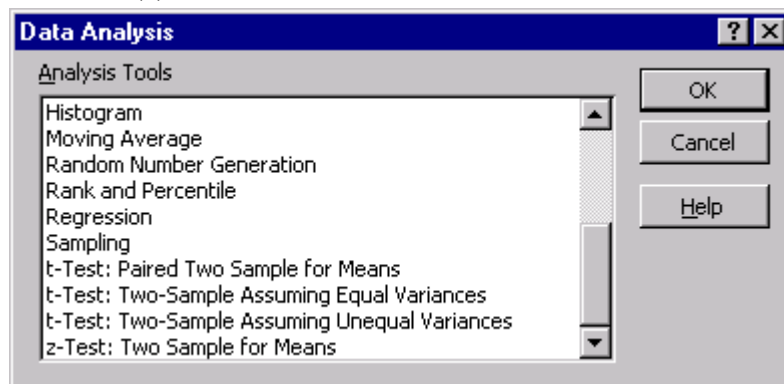
- a. Formulate hypotheses to test whether there is a significant difference between rates before and after the introduction of competition.

[5]

Let μ_1 = population average rate before, μ_2 = population average rate after

$H_0: \mu_1 = \mu_2$ $H_1: \mu_1 \neq \mu_2$.

- b. Indicate, on the following menu, which Excel Data Analysis Tool you would use to test the hypotheses in (a)



[5] t-Test: Paired Two Sample for Means

- c. Fill out the fields needed in this menu, to test the hypotheses in (a):

[5] B2:B5 C2:C5 Hypo Mean Diff = 0

- d. What assumptions are needed to use the above methods?

[5]

Random Sample, i.e. data pairs are independent.

Underlying population of differences are Normal.

- e. Here are two possible outputs from the appropriate Excel Data Analysis tool (one is right, the other is wrong):

t-Test: Paired Two Sample for Means			t-Test: Two-Sample Assuming Equal V:		
	Variable 1	Variable 2		Variable 1	Variable 2
Mean	22.0175	22.285	Mean	22.0175	22.285
Variance	7.941825	7.487233	Variance	7.941825	7.487233
Observations	4	4	Observations	4	4
df	3		df	6	
t Stat	-5.47220925		t Stat	-0.1362	
P(T<=t) one-tail	0.00599872		P(T<=t) one-tail	0.448058	
t Critical one-tail	2.35336302		t Critical one-tail	1.943181	
P(T<=t) two-tail	0.01199744		P(T<=t) two-tail	0.896116	
t Critical two-tail	3.18244929		t Critical two-tail	2.446914	

Choose the one that you think is right, and use it to give a p-value, with both yes-no ($\alpha = 0.01$) and gray level interpretations, to test the hypotheses in (a).

[10]

from Paired side: 2-sided p-value = 0.0120

yes-no: $\alpha = 0.01 > \text{p-val}$, so no strong evidence.

Gray level: evidence is strong, and quite close to very strong.

- f. Now focus only on the rates after the introduction of competition. Write an Excel formula for how large a sample size would be needed to estimate the mean of this population with an accuracy of 0.2, 90% of the time. (Hint: you can get needed information out of the tables in (e) above).

[10]

$$n = (\text{cutoff} * \text{sd} / B)^2 =$$

$$= (\text{NORMINV}(0.95,0,1) / 0.2)^2 * 7.94$$

- g. If Excel gives the numerical answer of 537.173407 to (f), what would you round it to?

[5]

round up to 538

3. 50% of all schools subscribe to CCN. Of these subscribers, 10% actually never use CCN, while 20% use CCN at least 5 times per week.

- a. Find the probability that a randomly selected school subscribes to CCN, but never uses it.

[5]

$$P\{\text{sub and never}\} = P\{\text{never} \mid \text{sub}\} P\{\text{sub}\} = (0.1)(0.5) = 0.05$$

- b. If a school is selected randomly from the subscribers, find the probability that it uses CCN less than 5 times per week.

[5]

$$P\{<5 \mid \text{sub}\} = 1 - P\{>=5 \mid \text{sub}\} = 1 - 0.2 = 0.8$$

- c. Find the probability that a randomly selected school never uses CCN (hint: this includes those that don't subscribe, and also those that subscribe, but don't use it).

[10]

$$\begin{aligned} P\{\text{never}\} &= P\{(\text{never and sub}) \text{ or } (\text{never and not sub})\} \\ &= P\{\text{never and sub}\} + P\{\text{never and not sub}\} \quad (\text{or rule}) \\ &= P\{\text{never} \mid \text{sub}\} P\{\text{sub}\} + P\{\text{not sub}\} \\ &= (0.1)(0.5) + 0.5 = 0.55 \end{aligned}$$

- d. If a randomly selected school doesn't use CCN, what is the probability they have subscribed to it?

[10]

$$\begin{aligned} P\{\text{sub} \mid \text{mean}\} &= P\{\text{sub and mean}\} / P\{\text{mean}\} \\ &= (0.1)(0.5) / ((0.1)(0.5) + 0.5) = 0.05 / 0.55 \end{aligned}$$

4. For the probability distribution:

x	0	1	2	3	4
$f(x)$	0.1	0.2	0.4	0.2	0.1

a. Find $P\{X \geq 1 | X < 3\}$.

[5]

$$\begin{aligned}
 &= P\{X \geq 1 \text{ and } X < 3\} / P\{X < 3\} \\
 &= P\{1 \leq X < 3\} / P\{X < 3\} \\
 &= (0.2 + 0.4) / (0.1 + 0.2 + 0.4) \\
 &= 0.6 / 0.7 \\
 &= 6/7
 \end{aligned}$$

b. Why is $EX = 2$?

[5]

$$\begin{aligned}
 EX &= \sum_x f(x) x = (0.1) 0 + (0.2) 1 + (0.4) 2 + (0.2) 3 + (0.1) 4 = \\
 &= 0 + 0.2 + 0.8 + 0.6 + 0.4 = 2
 \end{aligned}$$

or: notice “symmetric distribution around 2”.

c. Write down a calculation which shows that $\text{var}(X) = 1.2$.

[5]

$$\begin{aligned}
 \text{var}(X) &= E(X - EX)^2 = E(X - 2)^2 = \sum_x f(x) x^2 = \\
 &= (0.1) 2^2 + (0.2) 1^2 + (0.4) 0^2 + (0.2) 1^2 + (0.1) 2^2 \\
 &= 2(0.4 + 0.2) = 1.2
 \end{aligned}$$

d. What is the standard deviation of X ?

[5]

$$\text{sd}(X) = \sqrt{\text{var}(X)} = \sqrt{1.2}$$

e. For $g(x) = (x - 2)^2$, what is $Eg(X)$?

[5]

$$E(g(X)) = E(X - 2)^2 = 1.2, \quad \text{from part (c)}$$

5. To evaluate the accuracy of a Metlar scale, an item whose weight is known to be 14.01 ounces is weighed five times. The weights are entered in an Excel spreadsheet as shown here:

	H
37	14.04
38	14.01
39	13.99
40	14.03
41	14.02

- a. Write an Excel formula to calculate a p-value to test whether there is a statistically significant difference between the average value and 14.01.

[10]

Let μ = average weight. $H_0: \mu = 14.01$ $H_1: \mu \neq 14.01$

$$\begin{aligned} \text{P-value} &= P\{|Xbar - 14.01| = |average - 14.01 \text{ or m. c. } | B' \text{dry}\} \\ &= P\{|Xbar - 14.01| \geq |avg - 14.01| \mid \mu = 14.01\} \\ &= TDIST(ABS(AVERAGE(H37:H41)-14.01)/(STDEV(H37:H41)/SQRT(5)),4,2) \end{aligned}$$

- b. Interpret the result, if the answer to (a) is $p\text{-val} = 0.405$.

[5]

yes – no: no strong evidence
gray level: no strong evidence

- c. What assumptions are needed to work parts (a) and (b)?

[5]

Independent observations, i.e. a random sample
Individual measurements are normally distributed

- d. Write an Excel formula to give the endpoints of a 98% confidence interval for the mean.

[10]

$$\begin{aligned} &Xbar \pm \text{cutoff} * s / \text{sqrt}(n) \\ &= AVERAGE(H37:H31) - TINV(0.02,4) * STDEV(H37:H31)/SQRT(5) \\ &= AVERAGE(H37:H31) + TINV(0.02,4) * STDEV(H37:H31)/SQRT(5) \end{aligned}$$

6. a. A list of 101 numbers has $\bar{x} = 2$ and $s = 2$. Find $\sum_{i=1}^n x_i$ and $\sum_{i=1}^n x_i^2$.

[5]

$$\sum x_i = n \bar{x} = 101(2) = 202$$

$$4 = s^2 = (1 / (101 - 1)) [\sum x_i^2 - (101) 2^2]$$

$$= (1 / 100) [\sum x_i^2 - 404]$$

$$\text{so: } 400 = \sum x_i^2 - 404$$

$$\text{thus } \sum x_i^2 = 400 + 404 = 804$$

- b. Can a list of 10 numbers have $\sum_{i=1}^n x_i = 20$ and $\sum_{i=1}^n x_i^2 = 30$? Why or why not?

[5]

$$\bar{x} = \sum x_i / n = 20 / 10 = 2$$

$$\sum x_i^2 - n(\bar{x})^2 = 30 - 10(2)^2 = 30 - 40 < 0$$

so cannot have this.

7. In a random sample of 500 small business operators, 23% were motivated by a desire to be their own boss.

- a. Write down the steps you would use to check whether the sample size is large enough to use the Normal approximation for confidence intervals and hypothesis tests (but don't actually check).

[10]

$$n = 500, \quad \hat{p} = 0.23$$

$$\text{check that: } \hat{p} - 3 \sqrt{\hat{p}(1 - \hat{p}) / n} > 0$$

$$\text{and that: } \hat{p} + 3 \sqrt{\hat{p}(1 - \hat{p}) / n} < 1$$

- b. Suppose that the steps in (a) revealed that the Normal Approximation *is not* satisfactory. Write an Excel formula to calculate a p-value to determine whether the population percentage of all small business operators, who are motivated by a desire to be their own boss, is significantly more than 20%.

[5]

$$H_0: p \leq 0.2 \quad H_1: p > 0.2$$

$$\begin{aligned} \text{P-value} &= P\{\hat{p} = 0.23 \text{ or m.c.} \mid B; \text{dry}\} = P\{\hat{p} \geq 0.23 \mid p = 0.2\} \\ &= P\{X \geq n(0.23) \mid p = 0.2\} = 1 - P\{X \leq 500(0.23) - 1 \mid p = 0.2\} \\ &= 1 - \text{BINOMDIST}(500*0.23-1, 500, 0.2, \text{TRUE}) \end{aligned}$$

- c. Suppose that the steps in (a) revealed that the Normal Approximation *is* satisfactory. Repeat part (b), using a Normal approximation (with continuity correction).

[10]

$$= 1 - \text{NORMDIST}(500*0.23-0.5, 500*0.2, \text{SQRT}(500*0.2*0.8), \text{TRUE})$$

- d. Suppose that the steps in (a) revealed that the Normal Approximation *is* satisfactory. Write an Excel formula to give the endpoints of a level 90% conservative confidence interval for the population proportion of all small business operators, who are motivated by a desire to be their own boss.

[10]

$$\hat{p} \pm \text{cutoff} * \text{sqrt}(p*(1-p)) / \text{sqrt}(n)$$

For conservative: take $p = \frac{1}{2}$ (to max $p*(1-p)$)

$$= 0.23 - \text{CONFIDENCE}(0.05, \text{sqrt}(0.5*0.5), 500)$$

$$= 0.23 + \text{CONFIDENCE}(0.05, \text{sqrt}(0.5*0.5), 500)$$

- e. Suppose that the steps in (a) revealed that the Normal Approximation *is* satisfactory. Write an Excel formula to calculate how large (use the “best guess” method) a sample size should be used to estimate the population proportion so that the accuracy is within 0.005, with probability 0.98.

[5]

symmetric area inside is 0.98, when each tail has 0.01, so

$$= (\text{NORMINV}(0.01, 0, 1))^2 * 0.23 * 0.77 / 0.005^2$$