HW 3

If In a NY Times article (Feb 17, 1999) about the PSA blood test for detecting prostate cancer, of the men who had the disease, the test fails to detect prostate cancer in 1 in 4; of those without the cancer, 2/3 received false positives. Let C/C denote having/not having prostate cancer and let +/- denote a

positive/negative result.

a) Which is true: $P(-|C| = \frac{1}{4})$ or $P(C|-) = \frac{1}{4}$? $P(-|C| = \frac{1}{4})$ $P(\bar{C}|+) = \frac{2}{3} \text{ or } P(+|\bar{C}| = \frac{2}{3})$? $P(\bar{C}|+) = \frac{2}{3}$

b) What is the sensitivity of the test? Sensitivity = $P(+|C) = \frac{3}{4}$

c) of men who take the PCA test, suppose P(c) = .01. Find the cell probabilities in the $2 \times Z$ table for the joint distribution that cross classifies Y = diagnosis (+, -) with $X = true disease status (C, <math>\overline{C}$).

- d) Using (c), find the marginal distribution for the diagnosis. (above)
- e) Using (c) and (d), find P(C|+), and interpret P(C|+) = .6075 = 1

Approximately is men who test positive for prostate cancer will actually have the disease

Z For diagnostic testing, let X = true status (1= disease, Z = no disease) and Y = diagnosis (1= positive Z = negative). Let $\pi_i = P(Y = 1 | X = i)$, i = 1, Z

a) Explain why sensitivity = π , and specificity = $1-\pi_z$ sensitivity = π , because when Y=1, X=1 we have detected a true positive. specificity = $1-\pi_z$ because the complement of π_z is the case when true negatives are detected.

b) Let γ denote the probability the subject has the disease, Given that the diagnosis is positive, use Baye's theorem to show that the probability the subject truly has the disease is: $N, \gamma + N_2(1-\gamma)$

This is Bayes Theorem rewritten to correspond to our table format. $\pi, \pi_2 \pi, = P(C|T)$ and y = P(C)

 $|-\eta_1| |-\eta_2| | |\gamma_2| = P(\bar{c}|+) | |-\gamma| = P(\bar{c}|$

C) For mammograms for detecting breast cancer, suppose y = .01, sensitivity = .86, and specificity = .88. Given a positive test result, find the probability that a woman truly has breast cancer.

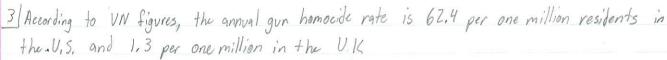
 $P(C) = (.86)(.01) = (.86)(.01) = (.86)(.01) = .009775 \approx .977\%$ $(.86)(.01) + (.88)(1 - .01) = (.86)(.01) = .009775 \approx .977\%$

d) To better understand the answer in (c), find the joint probabilities for the 2x2 cross classification of X and Y. Discuss their relative sizes in the two cells that refer to a positive result.

X + .86 .12 - .14 .88

The is much larger than Miz, which is good because we would not expect for a large number of "successes" to appear if the patient does not have the disease.

Otherwise our test would not be valuable because it would indicate that everyone has the disease, whether or not they do.



a) Compare the proportion of annual gun homocides using:

i) Difference of proportions

$$T_1 - T_2 = \frac{62.4}{1,000,000} - \frac{1.3}{1,000,000} = \frac{61.1}{1,000,000}$$

ii) Relative Risk

 $RR = \frac{T_1}{T_1c} = \frac{62.4}{1,000,000} = \frac{62.4}{1.3} = \frac{48}{1.3}$

b) When both proportions are very close to O, as here, which measurement is more useful for describing the strength of association? Why?

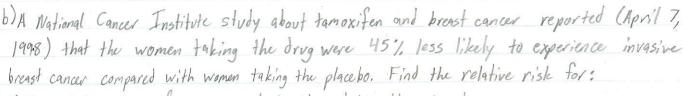
Relative Risk because (1) the resut is on a larger, more intuitive scale (everyday numbers instead of near zero ones), and (2) because it relates the proportions in terms of scale, rather than absolute difference.

4) A newspaper article preceding the 1994 World Cup semifinal match between Italy and Bulgaria stated that "Italy is favored 10-11 to beat Bulgaria, which is rated 10-3 to reach the final." Suppose this means that the odds Italy wins are 1/10 and that Bulgaria wins are 3/10. Find the probability each team wins and comment.

P(Italy wins) = 10 = 52381) These ortcomes are interesting because they would suggest that there is an outcome other than P(Bulgaria wins) = 30 = . 23077 either of these two countries winning the semifinal game. .52+.23 =1

5 | Consider the following two studies reported in the NY Times: a) A British study reported (Dec. 3, 1998) that, of smokers who get lung cancer, women were 1.7 times more vulnerable than men to get small-cell lung cancer." Is 1.7 a relative risk, or an odds ratio! A Relative Risk

* This is a relative risk because it compares two probabilities of two different events. An odds ratio compares the probability of success of one event to the probability of a failure of that event.



i) those taking tamoxifen compared to those taking the place to.

 $N_1 = (1-.45)N_2 \longrightarrow N_1 = (.55)$

ii) those taking the placebo and those taking tamoxifen.

6 In the U.S., the estimated annual probability that a woman over the age of 35 dies of lung cancer equals .001304 for current smokers and .000121 for non-smokers [M. Pagano and K. Gauvreau, Principles of Biostatistics, Belmont, CA: Duxbury Press (1993), p. 134]. a) Calculate and interpret the difference of proportions and the relative risk. Which is more informative for these data? Why?

DoP = (.001304) - (.000121) = .00118.

RR = (.001304) = 10.78

The RR is more informative because it is a more "real" number (more easily understood) than the DoP. It puts the information on a scale that is somewhat independent of the specific Probabilities.

b) Calculate and interpret the odds ratio. Explain why the RR and OR take similar values. odds smoker = (.001304) & .001306 odds non-smoker = (.000121) = .000121 1-(.000121) 1-(.001304)

Q= (.001306) (= 10.79 (151000.)

The OR and RR are similar because the probability of dying of cancer is near O, making the odds of each event very near the probability due to the near 1 divisor in the odds formula.

I For adults who sailed on the Titanic, the odds ratio between gender (male, female) and survival (yes, no) was 11.4.

a) What is wrong with the interpretation: The probability for survival for females was 11.4 times that for males"? Give the correct interpretation.

That interpretation would be for if the RR, the correct interpretation would be that the odds of a female to survive are 11.4 times the odds of a male surviving.

b) The odds of survival for females equaled 2.9. For each gender, find the proportion of survivars. OR = 1.4 odds male = odds finale = (2.9) = .254 OR (11.4) OR (11.4)

c) Find the value of R in the interpretation, "The probability of survival for females was R times that for males." $R = \frac{\pi_f}{\pi_m} = (.744) = 3.673$ (.203.)

8 A research study estimated that under a ceartain condition, the probability a subject would be referred for heart catheterization was . 906 for whites and . 847 for blacks.

a) A press release about the study stated the odds referral for cardiac catheterization are 60% of odds for whites. Explain how they obtained 60% (more accurately 57%).

This computation was done by correctly calculating the odds ratio

b) An AP story that described the study stated "Doctors were only 60% as likely to order cardiac catheterization for blacks as for whites." What is wrong with this interpretation? Give the correct percentage for this interpretation. (In stating results to the general public, it is better to use KK than UR, it is simpler to understand.)

This interpretation is using an OR in place of a RR. The true RR is:

RR = Phlack = (.847) = .935 -> 93.5%.

Publice (.906)

Joseph posted at the FBI website stated that of all blacks slain in 2005, 91% were slain by blacks, and of all whites slain in 2005, 83% were slain by whites. Let Y denote race of victim and X denote race of murderer.

Murdurer B .91 .17 W .83 .09 W .93 .09 B .17 ,91

a) Which conditional distribution do these statistics refer to, Ygiven X or Xgiven Y?

The Calculate the odds ratio between X and Y and interpret. $\hat{\theta} = \frac{(.91)(.83)}{(.09)(.17)} = 419.366$ This OR indicates that the odds of blacks being slain by blacks is almost 50 times that of whites being slain by blacks.

c) Given that the murderer was white, can you estimate the probability that the victim was white? What additional information would you need to do this? (Hint: How could you use Bayes' Theorem?).

No. If we were given the probability that the murderer was a white person, then using Bayes' therorem, using race of the murderer as the "test outcome (positive or negative) and race of the victim as the "disease" (white, not white) we could find the true probability that the victim was white.

12 A statistical analysis that combines information from several studies is called a meta analysis. A meta analysis compared aspirin and placebo on incidence of heart attack and of stroke, separately for men and for women. For the Women's Health Study, heart attacks were reported for 198 of 19,934 taking aspirin and for 193 of 19,942 taking placebo.

a) Construct the ZxZ table that classifies the treatment (aspirin, placebo) with

whether a heart attack was reported (yes, no).

	aspirin	placebo	
Sheart attack	198	193	391
sheart attack	19,736	19,749	39,485
	19,934	19,942	

b) Estimate the odds ratio. Interpret.

odds = $(\frac{198}{19,934}) = .01003$ ods = $(\frac{193}{19,749}) = .00987$ $\hat{\theta}_{1} = (.01003) = 1.01655$ $1 - (\frac{198}{19,934})$ $1 - (\frac{193}{19,749})$ (.00987)The odds of a woman having a heart attack appear to be slightly raised by taking aspirin.

c) Find a 95% confidence interval for the population odds ratio for women. Interpret. (As of 2006, results suggest that for women, aspirin was helpful for reducing risk of stroke, but not necessarily risk of heart attack.)

 $\ln \hat{\theta} \pm z_{\alpha_{2}}(SE) = \ln \left(1.01655\right) \pm \left(1.96\right) \frac{1}{1} + \frac{1}{1} + \frac{1}{1} + \frac{1}{1}$ $= \ln \left(1.01655\right) \pm \left(1.96\right) \frac{1}{198} + \frac{1}{19,934} + \frac{1}{193} + \frac{1}{199} + \frac{1}{199}$

 $= \ln(1.01655) \pm 1.96\sqrt{.01033} \rightarrow (-0.18279, 0.21562)$ $\rightarrow (e^{-.18279}, e^{-.21562}) = (.83294, 1.24063)$

It would not appear that it can be concluded that aspirin has any direct effect on whether or not a woman will have a heart attack.

13 Refer to table 2.1 about beleif in the afterlife. Table 2.1 Cross Classification of Beleif in Afterlife by Gender Beleif in Afterlife Gender Yes No or Undecided Total

F 509 116 625 n=1127

M 398 104 502 a) Construct a 90% confidence interval for the difference of proportions and interpret. $\hat{p_F} = \frac{509}{625} = .8144 \quad \hat{p_m} = \frac{398}{502} = .79283 \quad Dop = (.8144) - (.79283) = .02157$ CI: DoP ± Zaz SE = DoP ± 1.645 | P= (1-P=) + Pm(1-Pm) = (.02157) ± 1.645 (.8144)(.1856) + (.79283)(.207 nz $= .02157 \pm .03974 = (-.01767, .06081)$ It is inconclusive at the 90% confidence level whether women are more, or less likely to believe in an afterlife. b) Construct a 90% confidence interval for the odds ratio and interpret. odds_= .8144 = 4.38793 odds_= .79283 = 3.82695 \(\hat{\theta} = (4.38793) = 1.14659 1-.8144 (3.87695)CI: In (1,14659) + 1.645 / 1 + 1 + 1 + 1 = (-.11113, 38471) → (e, 1113) = (.89482, 1.46918) Again, we cannot conclude whether or not the odds are in favor of females beliving in the afterlife more often than men. c) Conduct a test of statistical independence. Report the P-Value and interpret. $\hat{\mathcal{U}} = nP_{i+} + p_{+} = n\left(\frac{n_{i+}}{n}\right)\left(\frac{n_{+}j}{n}\right) = \frac{n_{i+}n+j}{n}$ $\hat{\mathcal{U}}_{i,j} = (625)(907) = 502.995$ $\hat{n}_{12} = (625)(220) = 122.005$ $\hat{n}_{21} = (502)(907) = 404.005$ $\hat{n}_{22} = (502)(220) = 97.995$ (1177) $X^{2} = \sum \frac{(n_{ij} - u_{ij})^{2}}{u_{ij}} = \frac{(509 - 502.995)^{2} + (116 - 122.005)^{2} + (398 - 404.005)^{2} + (104 - 97.995)^{2}}{502.995} = .8245$ $df = (2-1)(2-1) = 1 \rightarrow |P>.250|$ This P-value suggests that the difference between males and females in beleit in afterlife is not significant