

# Statistics 2000, Section 001, Final (300 Points)

December 12, 2001, Dr. Jürgen Symanzik

Your Name: \_\_\_\_\_

First look at all 6 questions. Then start with the question that looks easiest to you. Continue with a more difficult question. Try to answer as many questions as possible in these 110 minutes.

Note that you will obtain at least partial credit if you indicate a correct formula but your final result is incorrect. If you just rely on your calculator without indicating the formula that should be used and your result is incorrect, you will obtain no credit at all for this part of a question.

## Question 1: Short Answers (60 Points)

1. Let  $z_i = \frac{x_i - \bar{x}}{s}$ ,  $i = 1, \dots, n$ , be the z-scores for  $n \geq 2$  arbitrary numbers  $x_1, x_2, \dots, x_n$  that are not all equal. Which of the following statements is correct: (10 Points)
  - (a) The correlation coefficient  $r$  between  $x$  and  $z$  must be negative since I subtract  $\bar{x}$  from each  $x_i$ .
  - (b) The correlation coefficient  $r$  between  $x$  and  $z$  is somewhere between -0.99 and -0.01.
  - (c) The correlation coefficient  $r$  between  $x$  and  $z$  is somewhere between 0.01 and 0.7.
  - (d) The correlation coefficient  $r$  between  $x$  and  $z$  is somewhere between 0.7 and 0.99.
  - (e) The correlation coefficient  $r$  between  $x$  and  $z$  is exactly +1.
  - (f) The correlation coefficient  $r$  between  $x$  and  $z$  is exactly 0.
  - (g) The correlation coefficient  $r$  between  $x$  and  $z$  is exactly -1.

2. Determine the slope and the  $y$ -intercept of the lines whose equations are given as:  
**(12 Points)**

(a)  $3x - 8y = 4$

Slope =

$y$ -intercept =

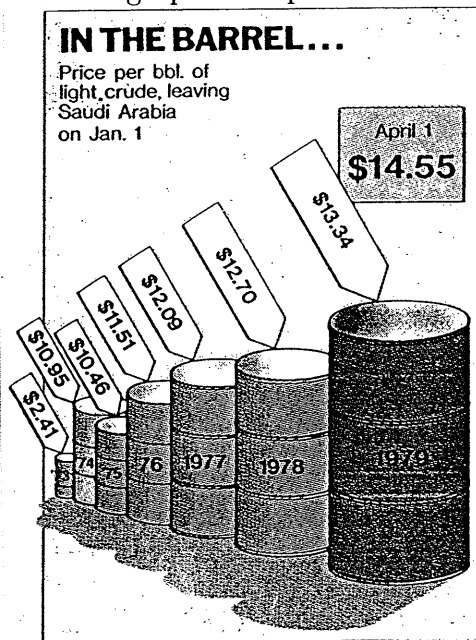
(b)  $4y + x + 5 = 0$

Slope =

$y$ -intercept =

3. A foreign lottery has a game called "4 out of 59". You make a selection of 4 numbers between 1 and 59 and you win the big prize if exactly these numbers are drawn in the weekly drawing. How many different combinations are possible to select 4 out of 59 numbers? Calculate this value! **(10 Points)**

- 6
4. What is wrong with this graphic from "Time" (April 9, 1979, p. 57), reprinted in Edward R. Tufte's book "The Visual Display of Quantitative Information". Provide a better graphical representation of the same data. (16 Points)



5. After Florida Secretary of State Katherine Harris certified the Florida results in the 2000 presidential election as a 537-vote Bush win, a Salt Lake City TV station held a call-in poll. The station asked its viewers to call in with their answer to the question:

"Do you believe Al Gore should stop trying to overturn legally certified votes in Florida and acknowledge that he has lost?"

4,237 people phoned in, and 3,482 said "yes". The station then announced that "82% of Americans believe that Al Gore should concede defeat."

Identify at least three problems with this survey and the announced results. (12 Points)

**Question 2: Normal Distribution (50 Points)**

**Part I:**

1. Let  $Z$  be a standard Normal variable, i.e.,  $Z \sim N(0, 1)$ , and  $X$  be a Normal variable with mean  $\mu = -4$  and variance  $\sigma^2 = 4$ , i.e.,  $X \sim N(-4, 2^2)$ . Determine the following: **(30 Points, i.e., 5 Points each)**

(a)  $P(Z < -2.2) =$

(b)  $P(-2.6 < Z < 1.6) =$

(c)  $P(X < -2.2) =$

(d)  $P(-2.6 < X < 1.6) =$

(e) Find a number # so that  
 $P(Z > \#) = 0.35$

(f) Find a number # so that  
 $P(X > \#) = 0.35$

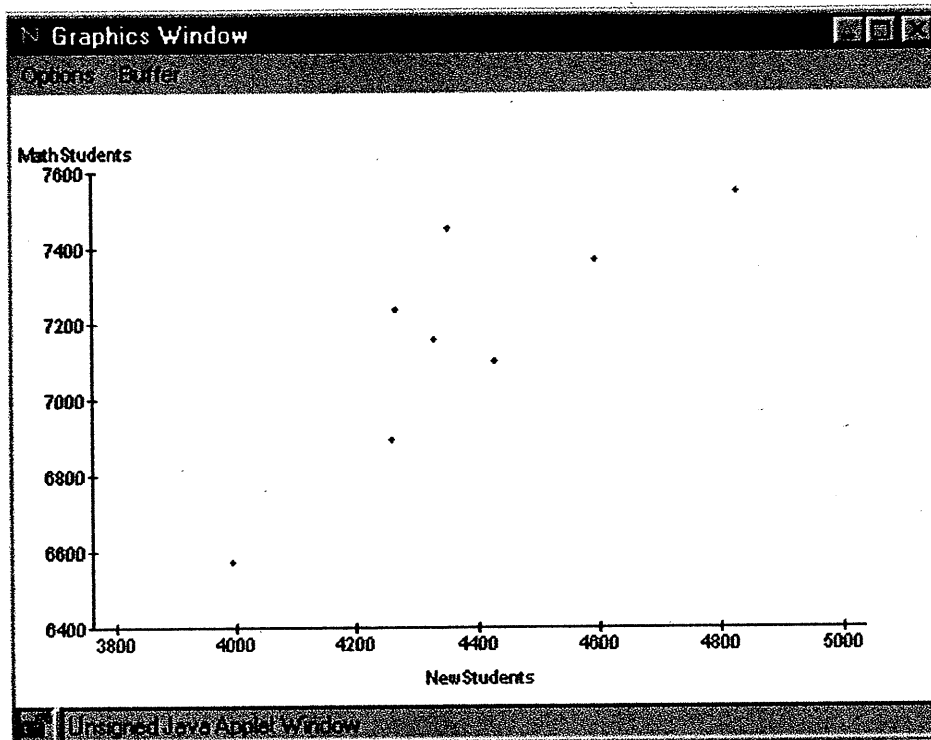


**Question 3: Linear Regression & Correlation (50 Points)**

The mathematics department of a large state university would like to use the number of freshman entering the university ( $\text{NewStudents} = x$ ) to predict the number of students who will sign up for freshman-level math courses ( $\text{MathStudents} = y$ ) in the fall semester. Data for the years 1986 through 1993 are given below as they appear in the WebStat main window. The fourth data column shows the Residuals related to the simple linear regression equation obtained on the next page.

N WebStat 2.0				
WebStat Data Stat Graphics Help				
	Year	NewStudents	MathStudents	Residuals
1	1986	4595	7364	-28.442984
2	1987	4827	7547	-92.82977
3	1988	4427	7099	-114.30083
4	1989	4258	6894	-139.09235
5	1990	3995	6572	-180.64957
6	1991	4330	7156	46.13244
7	1992	4265	7232	191.44339
8	1993	4351	7450	317.7397

The scatterplot of entering students (NewStudents) vs students taking freshman-level math courses (MathStudents) is shown below:



We used WebStat to fit a simple linear (least squares) regression to the data. Here is the numerical output:

Results:

Simple linear regression results:

Independent variable: NewStudents

Dependent variable: MathStudents

Sample size: 8

Correlation coefficient: 0.8333

(See fitted line plot in Graphics Panel.)

Residuals stored in column Residuals

Estimate of sigma: 188.94853

Parameter	Estimate	Std. Err.	DF	Tstat	Pval
Intercept	2492.6917	1267.1991	6	1.9670876	0.0484
NewStudents	1.0663223	0.2888466	6	3.691656	0.0051

1. Indicate the exact values for slope,  $y$ -intercept, the regression equation, and the correlation coefficient obtained from WebStat.

Provide an interpretation of Pearson's correlation coefficient  $r$  between entering students and students taking freshman-level math courses for this given data set. **(15 Points)**

2. Based on your equation in (1.), what is the predicted number of students taking freshman-level math courses when the number of entering students is 4,200?

And how many students taking freshman-level math courses would you predict when the number of entering students is 10,000? Explain! **(15 Points)**

3. Draw 2 different residual plots for this data set. One should contain the lurking variable. Is there any visible pattern in any of the your 2 residual plots? If there is a pattern, describe the pattern. Finally, argue whether the least squares regression line describes the relationship between entering students and students taking freshman-level math courses reasonably well. **(20 Points)**



3. Indicate  $P(Y \leq 3.0)$ ? (5 Points)

4. Calculate  $\mu = E(Y)$ , i.e., the mean (expected value) of  $Y$ . (10 Points)

5. Calculate  $\sigma^2 = Var(Y)$ , i.e., the variance of  $Y$ . (10 Points)

**Question 5: Binomial Probability Distributions (50 Points)**

In a previous quiz, 23 out of 25 Business Stat students answered a particular exam question correctly. Let us assume that in the upcoming academic year all new Business Stat students have to answer this same question in one of their quizzes. Based on previous years, the upper enrollment limit of 300 students has been obtained each year so we can safely assume that there will be exactly 300 students again during the next year that will attend this class.

Let  $X$  be the random variable that describes the number of students that will correctly answer this question in the next year. It is safe to assume that students' ability to answer a particular question does not change over time.

1. Indicate the probability distribution of  $X$ . Complete the following formula that relates to the probability distribution of  $X$ : **(5 Points)**

$$X \sim$$

2. Calculate the mean of  $X$  ( $\mu = E(X)$ ) and the variance of  $X$  ( $\sigma^2 = Var(X)$ ). **(15 Points)**

3. What is the **exact** probability that 290 students will answer this question correctly? Be as efficient in your calculations as possible. **(10 Points)**

4. What is the probability that 200 through 280 students will answer this question correctly? You may use an **approximation** if **appropriate** but need to check the required conditions first. **(20 Points)**

**Question 6: Tests of Significance (40 Points)**

1. Nutrition experts recommend that one's daily diet contain a minimum of 20 grams of fiber. The director of a summer camp for teenagers wants to show that the camp provides meals that exceed this amount. What null and alternative hypotheses should be tested? **(6 Points)**
  
2. Suppose a competitor believes that the camp's claim is an exaggeration, and that the camp is not meeting the nutritional needs of its participants with regard to fiber content. In particular, the competitor suspects that teenagers are provided a daily average of fewer than 20 grams of fiber. What set of hypotheses would the competitor be interested in testing? **(6 Points)**
  
3. Suppose an unbiased third party is only interested in determining if the camp's mean daily amount of fiber differs from 20 grams. It has no preconceived notion as to whether the actual mean is more or less than this figure. What set of hypotheses would this independent group want to test? **(6 Points)**
  
4. Suppose the unbiased third party is examining the fiber contents of meals on an annual basis. From previous years, it is known that the standard deviation is  $\sigma = 1.5$ . Based on a sample of 5 meals, the unbiased third party obtains a sample mean  $\bar{x} = 21.8$  grams of fiber.

Indicate the test statistic and calculate the  $p$ -value for the set of hypotheses specified in (3.). Can you reject  $H_0$  at the 5% level of significance? **(22 Points)**