

MATH 338

FINAL EXAM

THURSDAY, DECEMBER 15, 2016

Your name: _____

Your scores (to be filled in by Dr. Wynne):

Problem 1: ____/10

Problem 2: ____/19

Problem 3: ____/12

Problem 4: ____/9

Problem 5: ____/11

Problem 6: ____/12.5

Problem 7: ____/9

Total: ____/82.5

You have 110 minutes to complete this exam. This exam is closed book and closed notes with the exception of your two sheets of notes (front and back).

For full credit, show all work except for final numerical calculations (which can be done using a scientific calculator).

Problem 1. [1 pt each] Below are the names of a bunch of different hypothesis tests. For each claim in parts A-J, identify a correct test to use to test the claim (some claims may be tested using more than one test). Assume all assumptions for the tests are met. Tests may be used more than once or not at all.

- a. One sample t-test
- b. Two sample t-test
- c. Matched pairs t-test
- d. One sample proportion z-test
- e. Two sample proportion z-test
- f. Slope t-test for linear regression
- g. ANOVA test for linear regression
- h. One-way ANOVA
- i. Chi-square test for independence
- j. Chi-square test for goodness of fit

A. The distribution of Starburst colors is 25% red, 25% orange, 25% pink, and 25% yellow.

B. Gas costs more, per gallon, at Shell stations than at Mobil stations across the street from them.

C. There is no relationship between how much you eat (in calories) and your weight.

D. When people experience pain, they report less pain when they swear than when they don't swear.

E. One-third of American adults can't identify the United States on a world map.

F. Men are more likely to drive and text than women are.

G. Undergraduates in every CSUF College, on average, have the same GPA.

H. On average, a "gallon of milk" contains less than 1 gallon.

I. Knowing your favorite burger chain tells me nothing about whether you prefer pancakes or waffles.

J. The more money you have, on average, the more problems you have.

Problem 2. In a famous meta-analysis, Linus Pauling examined an experiment in which a group of children at ski school were given a pill containing either 1 gram of Vitamin C or a placebo, and were followed to see if they caught a cold. The data is shown in the two-way table below.

	Cold	No Cold	Total
Vitamin C	17	122	139
Placebo	31	109	140
Total	48	231	279

A. [1 pt] Why did the study use a placebo (instead of just not giving some children any pill at all)?

B. [2 pt] The study was a randomized and double-blind experiment. Explain what the two underlined terms mean in the context of this study.

Randomized:

Double-blind:

C. [7 pt] At the 5% significance level, is there a difference in the rate at which children catch colds when given Vitamin C compared to placebo?

Problem 2 (continued). The table is re-printed below for your convenience.

	Cold	No Cold	Total
Vitamin C	17	122	139
Placebo	31	109	140
Total	48	231	279

D. [2 pt] Under the null hypothesis for the chi-square test of independence, fill in the expected counts in the table below, to two decimal places:

	Cold	No Cold
Vitamin C		
Placebo		

E. [2 pt] Using the table you constructed in Part D, compute the chi-square test statistic.

In parts F-H, assume that one of the 279 children in the study is chosen at random.

F. [1 pt] What is the probability that the child was given a Vitamin C pill?

G. [2 pt] What is the probability that the child caught a cold, given that he or she got a Vitamin C pill?

H. [2 pt] What is the probability that the child got a Vitamin C pill, given that he or she caught a cold?

Problem 3. Colquhoun (2014) claimed that, “if you use $p=0.05$ to suggest that you have made a discovery, you will be wrong at least 30% of the time.” He uses the following example.

Suppose you pick a drug at random from 1000 candidate drugs and run a randomized controlled trial to determine whether the drug truly is effective. You perform a hypothesis test that has a significance level of 0.05 and a power of 0.8.

A. [2 pt] Write the null and alternative hypothesis for this hypothesis test, using words only (no math).

B. [1 pt] What is the probability of committing a Type I Error on this single hypothesis test?

C. [1 pt] What is the probability of committing a Type II Error on this single hypothesis test?

D. [3 pt] If you independently run experiments and hypothesis tests for four drugs, what is the probability of committing at least one Type I Error when performing the four hypothesis tests?

E. [5 pt] 10% of the drugs truly have an effect, and the other 90% are no different from placebo. Find the probability that your randomly chosen drug has an effect, given that you rejected the null hypothesis.

Problem 4. In 2009 the LA Times Data Desk started its “Mapping LA” project, which included collecting information about all 114 distinct neighborhoods of the city of Los Angeles. The dataset we will work with in Problems 4-7 includes the 104 neighborhoods that, as of 2008, had at least one public school in the neighborhood. Here is some important information about the variables in the dataset:

- Neighborhood: the name of the neighborhood
- Income: the median income of the neighborhood’s residents, in thousands of dollars
- Schools: the median API score of schools in the neighborhood, measured on a scale from 200 (bad) to 1000 (good)
- Asian, Black, Latino, White: the percentage of residents with each of those ethnicities (i.e., 5% Asian corresponds to Asian = 5)
- Diversity: how likely neighborhood residents are to encounter a resident of a different ethnicity, measured on a scale from 1 (not diverse at all) to 10 (very diverse)
- Population: the neighborhood’s population, according to the 2000 U.S. Census, in thousands

A. [1 pt] How many cases are in this dataset?

B. [1 pt] Which variable is the label variable?

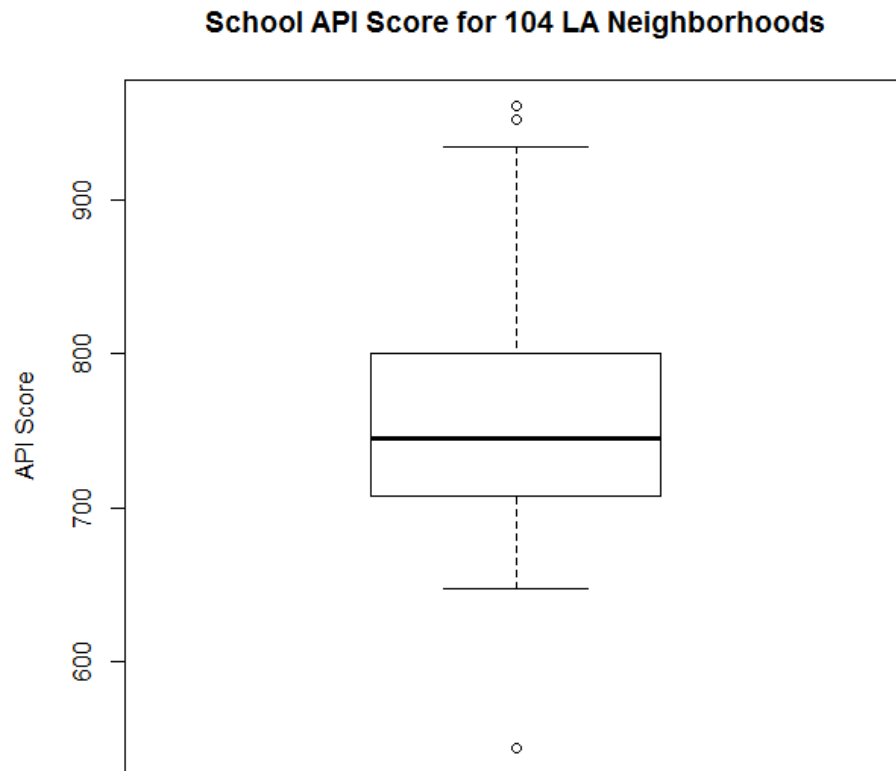
C. [1 pt] Is population a quantitative categorical variable (circle one)?

D. [1 pt] Was this dataset collected from an experiment observational study (circle one)?

E. [1 pt] Which of the following plots would be useful for identifying the distribution of Income:

 histogram pie chart scatter plot (circle one)?

Problem 4 (continued). The boxplot below plots the variable Schools. Use it to answer parts F-I.



F. [1 pt] Which of these values is closest to the median of the variable Schools:

400 650 700 750 800 (circle one)?

G. [1 pt] Which of these values is closest to the interquartile range of the variable Schools:

50 100 300 400 750 (circle one)?

H. [1 pt] How many outliers are there for the variable Schools?

I. [1 pt] Is the distribution of the variable Schools most likely:

skewed left skewed right symmetric (circle one)?

In Problems 5, 6, and 7, we will make some attempts at predicting how good a neighborhood's schools are (median API) from the other variables in the dataset.

Problem 5. Our first attempt is a multiple linear regression model using all five of the variables related to ethnicity (Asian, Black, Latino, White, and Diversity). Use the ANOVA table for the model (below) to answer parts A-D.

ANOVA Table

Source	DF	Sum of Squares	Mean Squares	F Value	Pr > F
Regression	5	404174	80834.7	41.9460	6.82533e-23
Residual	98	188857	1927.11		
Total	103	593031			

A. [2 pt] State the null and alternative hypothesis for the ANOVA F Test for this multiple linear regression model.

B. [1 pt] Find the F statistic and the p-value from the output above, and report their values below.

C. [1 pt] Based on your answer to part B, is the overall model significant at the 1% level?

D. [2 pt] Find the squared multiple correlation coefficient (R^2) for this multiple linear regression model. Round your answer to 2 decimal places.

Problem 5 (continued). Use the Parameter Estimates table for the model (below) to answer parts E-G.

Parameter Estimates

Variable	Parameter Estimate	Standard Error	t Value	Pr > t
(Intercept)	486.570	396.456	1.22730	0.222651
Diversity	-3.94718	1.98950	-1.98401	0.0500517
Asian	6.05296	4.08813	1.48062	0.141916
Black	2.04206	4.01382	0.508757	0.612066
Latino	1.93060	3.97734	0.485401	0.628476
White	3.90155	4.12266	0.946368	0.346289

E. [1 pt] If we performed backward selection, which variable would we remove from this model?

F. [2 pt] Interpret the value 6.05296 in the Parameter Estimates table.

G. [2 pt] Do you have any concerns about the choice of explanatory variables in this model? If so, what are they? If not, why not?

Problem 6. Our next attempt is a simple linear regression model using Income (in thousands of dollars) as the only predictor. Use the Parameter Estimates table (below) to answer all parts.

Parameter Estimates

Variable	Parameter Estimate	Standard Error	t Value	Pr > t
(Intercept)	655.641	11.9676	54.7847	1.81476e-77
Income	1.81123	0.188148	9.62660	5.46340e-16

- A. [1 pt] Write the equation used by this model to predict median school API from Income.
- B. [1.5 pt] Name three plots that we should look at to determine whether model assumptions are met.
- C. [3 pt] Construct (don't interpret) a 95% confidence interval for the population slope of the model.
- D. [2 pt] Predict the median school API for a neighborhood in which the median income is \$40,000.
- E. [5 pt] Income has a mean of 56.7312 and a variance of 835.483. The residual standard error (s) is 55.1934. Given these values and your answer to part D, construct and interpret a 90% prediction interval for the median school API in a neighborhood whose median income is \$40,000.

Problem 7. On our final attempt, we divide the neighborhoods into two categories: “mostly white” (White ≥ 50) and “mostly minority” (White < 50). Here is a summary of the data:

Group	Number of Neighborhoods	Mean(Schools)	sd(Schools)
Mostly Minority	70	725.53	54.60
Mostly White	34	826.06	68.77

A. [7 pt] Using a two-sample t test, test the claim that there is no difference in school performance between mostly minority and mostly white neighborhoods. Use a significance level of 1%.

B. [2 pt] Name one assumption necessary for the two-sample t test that is violated in this analysis. Justify why the assumption does not hold.

Extra Space. The tables below show a number of critical values z for the standard normal variable $Z \sim N(0, 1)$ and the corresponding cumulative proportions, corresponding to $P(Z \leq z)$.

z-score	Cumulative Proportion
-3.00	0.0013
-2.00	0.0228
-1.65	0.0495
-1.28	0.1003
-1.00	0.1587
-0.43	0.3336

z-score	Cumulative Proportion
0.43	0.6664
1.00	0.8413
1.28	0.8997
1.65	0.9505
2.00	0.9772
3.00	0.9987

Refer to the following tables for t^* and z^* critical values for confidence and prediction intervals:

Degrees of freedom	C = 0.90 (90%)	C = 0.95 (95%)	C = 0.98 (98%)	C = 0.99 (99%)
1	6.314	12.71	31.82	63.66
2	2.920	4.303	6.965	9.925
3	2.353	3.182	4.541	5.841
9	1.833	2.262	2.821	3.250
10	1.812	2.228	2.764	3.169
≈ 20	1.725	2.086	2.528	2.845
≈ 30	1.697	2.042	2.457	2.750
≈ 50	1.676	2.009	2.403	2.678
≈ 70	1.667	1.994	2.381	2.648
≈ 100	1.660	1.984	2.364	2.626
≈ 1000	1.646	1.962	2.330	2.581

	C = 0.90 (90%)	C = 0.95 (95%)	C = 0.98 (98%)	C = 0.99 (99%)
z^* values	1.645	1.960	2.326	2.576

For a two-sided hypothesis test, use the column corresponding to $C = 1 - \alpha$

For a one-sided hypothesis test, use the column corresponding to $C = 1 - 2\alpha$

Refer to the following table for χ^2 critical values:

Degrees of freedom	$\alpha = 0.05$	$\alpha = 0.01$
1	3.84	6.63
2	5.99	9.21
3	7.81	11.34
4	9.49	13.28