# Practice Midterm Exam

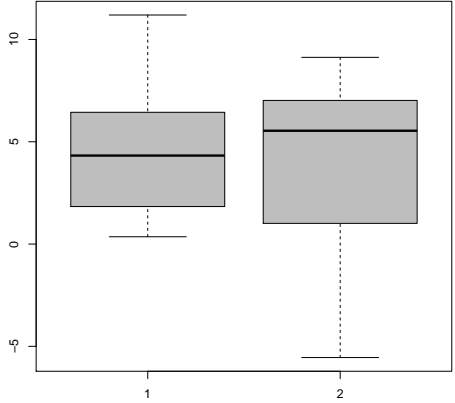### Statistics 139

Name _____ Harvard ID _____

This exam is closed book and closed notes, except that you may bring in and use two standard-sized sheets of paper (8.5" by 11") which can have notes on both sides. You are allowed to use a calculator. No copying, cheating, collaboration, or computers, or cell phones are allowed. The last 2 pages contain (1) a table of important distributions and (2) a table of probabilities from the standard normal distribution, and the page before that can be used for scratch work or for extra space. If you want any work done on the extra page or on backs of pages to be graded, mention where to look in big letters with a box around them, on the page with the question. Additionally:

1. Show your work and justify your answers. Explanations should be short, 2 or 3 sentences is often more than enough, but your reasoning should be clear and convincing.

2. Answers should be exact unless an approximation is asked for.

3. Simplify your answers unless otherwise specified, while keeping common sense in mind. For example, $\binom{52}{5}$ should be left as $\binom{52}{5}$ rather than doing a tedious calculation to get 2598960, but $\frac{\binom{52}{5}\cdot 2}{\binom{52}{5}\cdot 3}$ should be reduced to 2/3. Or you can use the numerical approximation of 0.667.

4. Use Type I error rates of $\alpha = 0.05$ and confidence levels of 95% unless explicitly stated otherwise. You can assume all tests are two-sided unless otherwise specified.

5. When performing a hypothesis test, make sure to (1) state the hypotheses, (2) state the calculated test statistic (and degrees of fredom if appropriate), (3) state the calculated p-value or critical value, and (4) state the conclusion in context of the problem along the the scope of inference.

## Good Luck!

| Problem | Maximum Points | Points Received |
|---------|----------------|-----------------|
| 1 | 15 | |
| 2 | 44 | |
| 3 | 16 | |
| 4 | 25 | |
| **Total** | **100** | |

**Problem 1. [3 points each] Parts are unrelated unless otherwise specified.**

(a) One scientist believes a new treatment will improve survival 6 months while another investigator believes it will improve 3 months over standard treatment. A valid $t$-based 95% confidence interval for a mean difference between 2 treatment and control was calculated to be (2, 8).

    A) Both scientistic claims are equally plausible since they are both inside the confidence interval.

    B) 3 months is more plausible since it is closer to the null hypothesis of a mean difference of 0.

    C) 6 months is more plausible since it is closer to the observed sample mean.

    D) Cannot be determined from the information given.

(b) You'd like to perform a hypothesis test to determine whether the the top 10% of earners in Massachusetts is different than it is in New Hampshire (as measured by the 90th percentile of income). A simple random sample of 500 income tax returns is taken within each state. Which test would make the most sense to perform to compare these two groups?

    A) Randomization test

    B) Permutation test

    C) Bootstrap test

    D) Proportion $z$-test

(c) You'd like to perform a test to determine whether the independent groups in the boxplot to the right come from distributions with similar centers. Which test is most appropriate?

    A) Unpooled $t$-test

    B) Permutation test

    C) Rank Sum Test

    D) ANOVA $F$-test

(d) A test of $H_0 : \mu = 0$ versus $H_A : \mu \neq 0$ is conducted on the same population independently by two different researchers. They both use the same sample size and the same value of $\alpha = 0.05$. Which of the following will be the same for both researchers?

    A) The $p$-value of the test.

    B) The power of the test if the true $\mu = 6$.

    C) The value of the test statistic.

    D) The decision about whether or not to reject the null hypothesis.

(e) In a hypothesis test the decision was made to not reject the null hypothesis. Which type of mistake could have been made?

    A) Type 1.

    B) Type 2.

    C) Type 1 if it's a one-sided test and Type 2 if it's a two-sided test.

**Problem 2. [5 points each unless stated] Parts are unrelated unless otherwise stated.**

(a) A recent study claimed that getting rid of annual medical check-ups would be harmful since patients that attend their annual check-ups have better health outcomes than those patients that skip them. Provide one reason why this claim may be incorrect. Be specific.
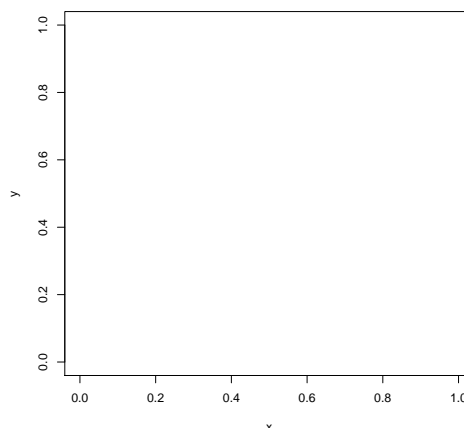
(b) You perform a hypothesis test of the mean using a sample size of four units, and you do not reject the null hypothesis. Your research colleague says this statistical test provides conclusive evidence that the null hypothesis is reasonable. Do you agree or disagree with his conclusion? Explain your reasoning in three or fewer sentence.

(c) In comparing 10 groups, you notice that $Y_7$ is the largest and $Y_3$ is the smallest, and proceed to test the hypothesis that $H_0 : \mu_3 - \mu_7 = 0$. Why should a multiple comparison procedure be used even though there is only one comparison being made?

(d) The median test score on a Stat 139 exam was 85 (out of 100). Would the mean be expected to be above, below, or around the same value of 85? In one or two sentences, explain why.

(e) In the plot provided below, draw a scatter of points that clearly violates one of the assumptions of linear regression, but not the others. Be sure to mention which assumption it violates.

Assumption Violated:

_____



(f) A regression is run in order to determine whether the last $n = 102$ monthly returns of Microsoft Stock prices (`msft`) mimic that of McDonald's Stock prices (`mcd`). Based on the R-output below, determine whether a slope of $\beta_1 = 1$ is reasonable.

```
> summary(lm(msft~mcd))

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.000278   0.006720  -0.041    0.967
mcd          0.783439   0.160266   4.888 3.97e-06
```

(g) Let $t = \dfrac{\hat{\beta}_1}{\hat{\sigma}_e \sqrt{\frac{1}{(n-1)S_X^2}}}$ be the usual $t$-statistic for the OLS slope estimate in a simple regression model. Derive the distribution of $t^2$ through representation. What assumptions do you need for this result to be exact?

(h) In the multiple regression model based on OLS with no missing data (briefly explain your answer):

i. True or False: $R^2$ cannot be negative.

ii. True or False: $R^2$ is the percent of variability in the response variable associated with the predictors.

iii. True or False: $R^2$ equals the sum of the squares of the separate correlation coefficients $r$ of the response with each predictor separately.

iv. True or False: $R^2$ may decrease when an additional explanatory variable is added.

4

**Problem 3. [16 points total]**

A friend tells you that she is the master of the game 'Rock-Paper-Scissors' and reports that she truly wins each round of the game with 0.8 probability, but you think she is just bragging (or even lying) and really wins with 0.5 probability. You decide to test your friend on her claim by playing 25 independent rounds of the game, and decide to believe her if she wins 18 or more of the rounds.

(i) Write down the hypotheses, the test statistic, and determine the true reference distribution for this test statistic for this study.

(ii) Calculate an approximate Type 1 error rate for this test.

(iii) Calculate an approximate power for this test.

(iv) You play your 25 rounds of the game and she wins 16 of them. What do you conclude statistically from this test? What do you conclude practically?

**Problem 4. [25 points total]** The following is the R-output for a regression to predict the number of bowls of noodle soup sold at a hip new pho restaurant in town based on the high temperature outside the restuarant that day:

```
> mean(temp)
[1] 52.7
> sd(temp)
[1] 17.35
> mean(soup)
[1] 225.12
> sd(soup)
[1] 86.82
> summary(lm(soup~temp))

Estimate Std. Error t value Pr(>|t|)
(Intercept)    ------      ------   26.84   <2e-16 ***
  temp              ------      ------  -14.51   <2e-16 ***

  Residual standard error: 37.80 on 48 degrees of freedom
Multiple R-squared:  0.8143,    Adjusted R-squared:  0.8105
F-statistic: 210.5 on 1 and 48 DF,  p-value: < 2.2e-16
```
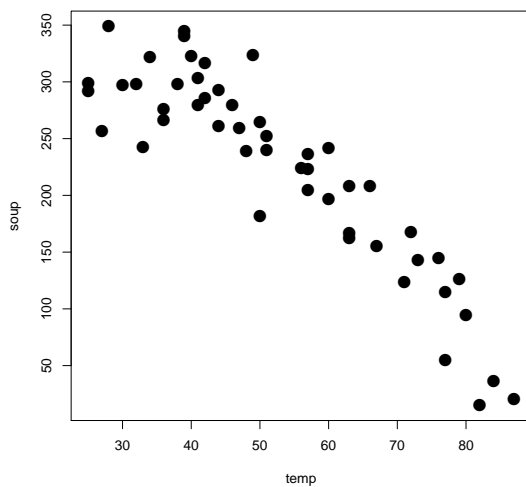
(a) [3 points] What is the estimated correlation between `temp` and `soup`?

(b) [4 points] Calculate the estimated simple regression line for these data.

(c) [5 points] Calculate a 95% confidence interval for estimating the mean number of bowls of soup sold when the high temperature outside is 52.7 degrees.

(d) [4 points] For a new randomly sampled day when the high temperature is known to be $X = 52.7$, what is the approximate probability that the number of bowls sold will be within the interval in part (c), assuming all assumptions are correct?

(e) [3 points] Here's the scatterplot of these data:



Provide the best set of transformations on $Y = soup$ and $X = temp$ to make a better regression model. Explain your choice.

(f) [6 points] A second predictor, `heat_index`, is considered to be used to predict soup sales. If `heat_index` is equal to `10*(temp - 90)`, then:

i. What are the estimates of the intercept and slope if `heat_index` is used as the only predictor for `soup`

ii. Give an estimate of the value of $R^2$ for a multiple regression model to predict `soupd` from both `heat_index` and `temp`? Explain how you came to that estimate.

Table of Distributions

| Name | Paramater | PMF or PDF | Mean | Variance |
|---|---|---|---|---|
| Bernoulli | $p$ | $P(X=1)=p,$ <br> $P(X=0)=1\text{-}p$ | $p$ | $p(1\text{-}p)$ |
| Binomial | $n, p$ | $\binom{n}{k}p^k(1-p)^{n-k}; k \in \{0,1,....n\}$ | $np$ | $np(1\text{-}p)$ |
| Geometric | $p$ | $p(1-p)^k; k \in \{0,1,2,...\}$ | $(1\text{-}p)/p$ | $(1\text{-}p)/p^2$ |
| Negative Binomial | $r, p$ | $\binom{r+k-1}{r-1}p^r(1-p)^k; k \in \{0,1,....\}$ | $r(1\text{-}p)/p$ | $r(1\text{-}p)/p^2$ |
| Hyper-geometric | $w, b, n$ | $\binom{w}{k}\binom{b}{n-k}\Big/\binom{w+b}{n}; k \in \{0,1,....n\}$ | $\mu = \dfrac{nw}{w+b}$ | $\left(\dfrac{w+b-n}{w+b-1}\right)n\left(\dfrac{\mu}{n}\right)\left(1-\dfrac{\mu}{n}\right)$ |
| Multinomial | $p_1,\ldots,p_k$ <br> $\sum p_i = 1$ | $\dfrac{n!}{x_1!x_2!...x_k!}p_1^{x_1}p_2^{x_2}...p_k^{x_k}$ | $E(X_i)=np_i$ | $Var(X_i)=np_i(1-p_i)$ <br> $Cov(X_i,X_j)=-np_i\,p_j,$ <br> for $i \neq j$ |
| Poisson | $\lambda$ | $\dfrac{e^{-\lambda}\lambda^k}{k!}; k \in \{0,1,....\}$ | $\lambda$ | $\lambda$ |
| Uniform | $a < b$ | $\dfrac{1}{b-a}; a < x < b$ | $(a+b)/2$ | $(b\text{-}a)^2/12$ |
| Normal | $\mu, \sigma^2$ | $\dfrac{1}{\sigma\sqrt{2\pi}}\exp\left(-(x-\mu)^2/(2\sigma^2)\right)$ | $\mu$ | $\sigma^2$ |
| Log-Normal | $\mu, \sigma^2$ | $\dfrac{1}{x\sigma\sqrt{2\pi}}\exp\left(-(\log(x)-\mu)^2/(2\sigma^2)\right)$ <br> $; x > 0$ | $\theta=\exp(\mu+\sigma^2/2)$ | $\theta^2[\exp(\sigma^2)\text{-}1]$ |
| Exponential | $\lambda$ | $\lambda\exp(-\lambda x); x > 0$ | $1/\lambda$ | $1/\lambda^2$ |
| Gamma | $a, \lambda$ | $\dfrac{\lambda^a}{\Gamma(a)}x^{a-1}\exp(-\lambda x)$ | $a/\lambda$ | $a/\lambda^2$ |
| Beta | $a, b$ | $\dfrac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}x^{a-1}(1-x)^{b-1}; 0 < x < 1$ | $\mu = \dfrac{a}{a+b}$ | $\dfrac{\mu(1-\mu)}{a+b+1}$ |
| Chi-Square | $d$ | $\dfrac{1}{2^{d/2}\Gamma(d/2)}x^{(d/2)-1}\exp(-x/2);$ <br> $x > 0$ | $d$ | $2d$ |
| $t$ | $d$ | $\dfrac{\Gamma((d+1)/2)}{\sqrt{d\pi}\Gamma(d/2)}(1+x^2/d)^{-(d+1)/2}$ | 0 if $d > 1$ | $d/(d\text{-}2)$ if $d > 2$ |
| $F$ | $d_1, d_2$ | $\dfrac{\Gamma(d_1+d_2)}{\Gamma(d_1)\Gamma(d_2)x}\sqrt{\dfrac{(d_1x)^{d_1}d_2^{d_2}}{(d_1x+d_2)^{d_1+d_2}}}$ | $d_2/(d_2\text{-}2)$ <br> if $d_2 > 2$ | $2d_2^2(d_1+d_2\text{-}2)/[d_1(d_2\text{-}2)^2(d_2\text{-}4)]$ <br> if $d_2 > 4$ |

## Useful R output:

```r
x=seq(0,3,0.01)
p=round(pnorm(x),4)
names(p)=x
p
```

```
##      0    0.01   0.02   0.03   0.04   0.05   0.06   0.07   0.08   0.09    0.1   0.11   0.12   0.13
## 0.5000 0.5040 0.5080 0.5120 0.5160 0.5199 0.5239 0.5279 0.5319 0.5359 0.5398 0.5438 0.5478 0.5517
##   0.14   0.15   0.16   0.17   0.18   0.19    0.2   0.21   0.22   0.23   0.24   0.25   0.26   0.27
## 0.5557 0.5596 0.5636 0.5675 0.5714 0.5753 0.5793 0.5832 0.5871 0.5910 0.5948 0.5987 0.6026 0.6064
##   0.28   0.29    0.3   0.31   0.32   0.33   0.34   0.35   0.36   0.37   0.38   0.39    0.4   0.41
## 0.6103 0.6141 0.6179 0.6217 0.6255 0.6293 0.6331 0.6368 0.6406 0.6443 0.6480 0.6517 0.6554 0.6591
##   0.42   0.43   0.44   0.45   0.46   0.47   0.48   0.49    0.5   0.51   0.52   0.53   0.54   0.55
## 0.6628 0.6664 0.6700 0.6736 0.6772 0.6808 0.6844 0.6879 0.6915 0.6950 0.6985 0.7019 0.7054 0.7088
##   0.56   0.57   0.58   0.59    0.6   0.61   0.62   0.63   0.64   0.65   0.66   0.67   0.68   0.69
## 0.7123 0.7157 0.7190 0.7224 0.7257 0.7291 0.7324 0.7357 0.7389 0.7422 0.7454 0.7486 0.7517 0.7549
##    0.7   0.71   0.72   0.73   0.74   0.75   0.76   0.77   0.78   0.79    0.8   0.81   0.82   0.83
## 0.7580 0.7611 0.7642 0.7673 0.7704 0.7734 0.7764 0.7794 0.7823 0.7852 0.7881 0.7910 0.7939 0.7967
##   0.84   0.85   0.86   0.87   0.88   0.89    0.9   0.91   0.92   0.93   0.94   0.95   0.96   0.97
## 0.7995 0.8023 0.8051 0.8078 0.8106 0.8133 0.8159 0.8186 0.8212 0.8238 0.8264 0.8289 0.8315 0.8340
##   0.98   0.99      1   1.01   1.02   1.03   1.04   1.05   1.06   1.07   1.08   1.09    1.1   1.11
## 0.8365 0.8389 0.8413 0.8438 0.8461 0.8485 0.8508 0.8531 0.8554 0.8577 0.8599 0.8621 0.8643 0.8665
##   1.12   1.13   1.14   1.15   1.16   1.17   1.18   1.19    1.2   1.21   1.22   1.23   1.24   1.25
## 0.8686 0.8708 0.8729 0.8749 0.8770 0.8790 0.8810 0.8830 0.8849 0.8869 0.8888 0.8907 0.8925 0.8944
##   1.26   1.27   1.28   1.29    1.3   1.31   1.32   1.33   1.34   1.35   1.36   1.37   1.38   1.39
## 0.8962 0.8980 0.8997 0.9015 0.9032 0.9049 0.9066 0.9082 0.9099 0.9115 0.9131 0.9147 0.9162 0.9177
##    1.4   1.41   1.42   1.43   1.44   1.45   1.46   1.47   1.48   1.49    1.5   1.51   1.52   1.53
## 0.9192 0.9207 0.9222 0.9236 0.9251 0.9265 0.9279 0.9292 0.9306 0.9319 0.9332 0.9345 0.9357 0.9370
##   1.54   1.55   1.56   1.57   1.58   1.59    1.6   1.61   1.62   1.63   1.64   1.65   1.66   1.67
## 0.9382 0.9394 0.9406 0.9418 0.9429 0.9441 0.9452 0.9463 0.9474 0.9484 0.9495 0.9505 0.9515 0.9525
##   1.68   1.69    1.7   1.71   1.72   1.73   1.74   1.75   1.76   1.77   1.78   1.79    1.8   1.81
## 0.9535 0.9545 0.9554 0.9564 0.9573 0.9582 0.9591 0.9599 0.9608 0.9616 0.9625 0.9633 0.9641 0.9649
##   1.82   1.83   1.84   1.85   1.86   1.87   1.88   1.89    1.9   1.91   1.92   1.93   1.94   1.95
## 0.9656 0.9664 0.9671 0.9678 0.9686 0.9693 0.9699 0.9706 0.9713 0.9719 0.9726 0.9732 0.9738 0.9744
##   1.96   1.97   1.98   1.99      2   2.01   2.02   2.03   2.04   2.05   2.06   2.07   2.08   2.09
## 0.9750 0.9756 0.9761 0.9767 0.9772 0.9778 0.9783 0.9788 0.9793 0.9798 0.9803 0.9808 0.9812 0.9817
##    2.1   2.11   2.12   2.13   2.14   2.15   2.16   2.17   2.18   2.19    2.2   2.21   2.22   2.23
## 0.9821 0.9826 0.9830 0.9834 0.9838 0.9842 0.9846 0.9850 0.9854 0.9857 0.9861 0.9864 0.9868 0.9871
##   2.24   2.25   2.26   2.27   2.28   2.29    2.3   2.31   2.32   2.33   2.34   2.35   2.36   2.37
## 0.9875 0.9878 0.9881 0.9884 0.9887 0.9890 0.9893 0.9896 0.9898 0.9901 0.9904 0.9906 0.9909 0.9911
##   2.38   2.39    2.4   2.41   2.42   2.43   2.44   2.45   2.46   2.47   2.48   2.49    2.5   2.51
## 0.9913 0.9916 0.9918 0.9920 0.9922 0.9925 0.9927 0.9929 0.9931 0.9932 0.9934 0.9936 0.9938 0.9940
##   2.52   2.53   2.54   2.55   2.56   2.57   2.58   2.59    2.6   2.61   2.62   2.63   2.64   2.65
## 0.9941 0.9943 0.9945 0.9946 0.9948 0.9949 0.9951 0.9952 0.9953 0.9955 0.9956 0.9957 0.9959 0.9960
##   2.66   2.67   2.68   2.69    2.7   2.71   2.72   2.73   2.74   2.75   2.76   2.77   2.78   2.79
## 0.9961 0.9962 0.9963 0.9964 0.9965 0.9966 0.9967 0.9968 0.9969 0.9970 0.9971 0.9972 0.9973 0.9974
##    2.8   2.81   2.82   2.83   2.84   2.85   2.86   2.87   2.88   2.89    2.9   2.91   2.92   2.93
## 0.9974 0.9975 0.9976 0.9977 0.9977 0.9978 0.9979 0.9979 0.9980 0.9981 0.9981 0.9982 0.9982 0.9983
##   2.94   2.95   2.96   2.97   2.98   2.99      3
## 0.9984 0.9984 0.9985 0.9985 0.9986 0.9986 0.9987
```

**Scratch Work:**