

Quiz 1, STATS 401 W18

In lab on 10/5

This document produces different random quizzes each time the source code generating it is run. The actual quiz will be a realization generated by this random process, or something similar.

This version lists all the questions currently in the quiz generator. The quiz will have one question drawn at random from each of the five categories. No new questions will be added after Wednesday 10/3. Small changes may be made.

Instructions. You have a time allowance of 40 minutes, though the quiz may take you less time and you can leave lab once you are done. The quiz is closed book, and you are not allowed access to any notes. Any electronic devices in your possession must be turned off and remain in a bag on the floor.

Formulas

The following formulas will be provided. To use these formulas properly, you need to make appropriate definitions of the necessary quantities.

(1) $\mathbf{b} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{y}$

(2) $\text{Var}(X) = E[(X - E[X])^2] = E[X^2] - (E[X])^2$

(3) The probability density function of the standard normal distribution is $\frac{1}{\sqrt{2\pi}} e^{-x^2/2}$

(4) Syntax from `?pnorm`:

```
pnorm(q, mean = 0, sd = 1)
qnorm(p, mean = 0, sd = 1)
q: vector of quantiles.
p: vector of probabilities.
```

(5) If a random variable is normally distributed, the probability it falls within one standard deviation of the mean is 68%, within two standard deviations of the mean is 95%, and within three standard deviations of the mean is 99.7%.

Q1. Matrix exercises

Q1-1.

(a). Evaluate $\mathbb{A}\mathbb{B}$ when

$$\mathbb{A} = \begin{bmatrix} 2 & 3 \\ 1 & 3 \\ -1 & -2 \end{bmatrix}, \quad \mathbb{B} = \begin{bmatrix} 3 & 1 \\ 0 & 1 \end{bmatrix}$$

(b). For \mathbb{A} as above, write down \mathbb{A}^T .

(c). For \mathbb{B} as above, find \mathbb{B}^{-1} if it exists. If \mathbb{B}^{-1} doesn't exist, explain how you know this.

Q1-2.

(a). Evaluate $\mathbb{A}\mathbb{B}$ when

$$\mathbb{A} = \begin{bmatrix} -1 & -1 & 3 \\ 2 & 0 & 3 \end{bmatrix}, \quad \mathbb{B} = \begin{bmatrix} -1 & 1 & -2 \\ 0 & 0 & 0 \\ -2 & 3 & 0 \end{bmatrix}$$

(b). For \mathbb{A} as above, write down \mathbb{A}^T .

(c). For \mathbb{A} as above, find \mathbb{A}^{-1} if it exists. If \mathbb{A}^{-1} doesn't exist, explain how you know this.

Q2. Summation exercises

Q2-1.

Calculate $\sum_{i=k}^{k+3} (i+3)$, where k is a whole number. Your answer should depend on k .

Q2-2.

Evaluate $\sum_{i=1}^{30} 10 - \sum_{i=10}^{20} 20$.

Q2-3.

Calculate $\sum_{k=m}^n a$, where m and n are whole numbers and a is a real number.

Q2-4.

Evaluate $3 \sum_{k=1}^5 2 - 0.5 \sum_{i=2}^{11} 6$.

Q2-5.

Evaluate $\sum_{i=1}^3 i(i-1)$.

Q2-6.

Suppose $F_0 = 0, F_1 = 1, F_2 = 1, F_3 = 2, F_4 = 3, F_5 = 5, F_6 = 8, F_7 = 13$ Evaluate $\sum_{i=4}^7 F_i - \sum_{i=0}^3 F_i$.

Q3. R exercises

Q3-1.

(a) Which of the following is the output of `matrix(c(rep(0,times=4),rep(1,times=4)),ncol=2)`

$$\begin{array}{llll} \text{(i). } \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} & \text{(ii). } \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} & \text{(iii). } \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix} & \text{(iv). } \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix} \end{array}$$

(b) Suppose we define an R vector by `y <- c(3,NA,-1,4,NA,-2)`. What will `y[y>0]` give you?

- (i). A vector of the positive elements and NA values of `y`.
 - (ii). A vector of the negative elements of `y`.
 - (iii). A vector of all NAs.
 - (iv). A vector of TRUEs and FALSEs.
 - (v). A vector of TRUEs and FALSEs and NAs.
-

Q3-2.

(a) Which one of the following lines of code successfully constructs the matrix $\mathbb{A} = \begin{bmatrix} 1 & 1 \\ 2 & 2 \\ 3 & 3 \end{bmatrix}$

- (i). `A <- matrix(c(1,1,2,2,3,3) ,nrow=3)`
- (ii). `A <- cbind(c(1,1),c(2,2),c(3,3))`
- (iii). `A <- t(matrix(c(1,1,2,2,3,3) ,nrow=2))`
- (iv). `A <- c(c(1:3),c(1:3))`

(b) Suppose `X` is a matrix in R. Which of the following is NOT equivalent to `X`?

- (i). `t(t(X))`
 - (ii). `X %*% matrix(1,ncol(X))`
 - (iii). `X*1`
 - (iv). `X%*%diag(ncol(X))`
-

Q3-3.

(a) Which of the following is the matrix `A` generated by

```
A <- t(matrix(c(rep(1,times=2),rep(3,times=2), 6, 4),ncol=3))
```

$$(i) \quad \mathbb{A} = \begin{bmatrix} 1 & 1 \\ 3 & 3 \\ 6 & 4 \end{bmatrix}$$

$$(ii) \quad \mathbb{A} = \begin{bmatrix} 1 & 3 & 6 \\ 1 & 3 & 4 \end{bmatrix}$$

$$(iii) \quad \mathbb{A} = \begin{bmatrix} 1 & 3 \\ 1 & 6 \\ 1 & 3 \end{bmatrix}$$

$$(iv) \quad \mathbb{A} = \begin{bmatrix} 1 & 1 & 3 \\ 3 & 6 & 4 \end{bmatrix}$$

(b) Which of the following successfully select the first five odd elements of the vector `x <- c(1,2,3,4,5,6,7,8,9,10,11)`? (List all that apply. Do not list commands that will give an error)

- (i) `x[rep(c(TRUE,FALSE),each=5)]`
 - (ii) `x[rep(c(TRUE,FALSE),times=5)]`
 - (iii) `x[rep(c(TRUE,FALSE),length=9)]`
 - (iv) `x[rep(c(TRUE,FALSE))][1:5]`
 - (v) `x[rep(c("TRUE","FALSE"),5)]`
 - (vi) None of the above
 - (vii) All of the above
-

Q3-4.

(a) Define the matrix `A` as:

```
##      [,1] [,2]
## [1,]    0    3
## [2,]    1    3
## [3,]    1    2
```

What is the output of `apply(A,2,mean)`?

- (i). A vector of length 3 corresponding to the average of each row of `A`.

- (ii). A vector of length 2 corresponding to the average of each column of **A**.
 - (iii). The mean of all the values in **A**.
 - (iv). The mean of the second column of **A**.
 - (v). The mean of the second row of **A**.
- (b) For each of the lines of code below, say whether it will correctly make 50 draws from the normal(100, 20) distribution. Among the correct answers, comment briefly on some strengths and weaknesses from the perspective of writing good R code. Which answer do you think is the best code, and why?

- (i) `rnorm(50,20,100)`
 - (ii) `rnorm(100,20,50)`
 - (iii) `rnorm(100,20,n=50)`
 - (iv) `rnorm(mean=100,sd=20,n=50)`
 - (v) `rnorm(n=50,mean=100,sd=20)`
 - (vi) `replicate(rnorm(100,20),50)`
 - (vii) `replicate(rnorm(n=1,mean=100,sd=20),n=50)`
 - (viii) `rnorm(50)*20+100`
 - (ix) `100+sqrt(20)*rnorm(50)`
-

Q3-5.

- (a) Which of the following successfully select the diagonal elements of the matrix

$A = \begin{bmatrix} 1 & 0 \\ 2 & 2 \end{bmatrix}$ represented in R by `A<-matrix(c(1,2,0,2),2,2)`?

- (i). `A[c(1,1),c(2,2)]`
- (ii). `A[rbind(c(1,1),c(2,2))]`
- (iii). `A[cbind(c(1,1),c(2,2))]`
- (iv). `A[matrix(c(TRUE,FALSE,FALSE,TRUE),2)]`
- (v). all of (i,ii,iii,iv)

- (vi). none of (i,ii,iii,iv)
 - (vii). (ii) and (iv) only
 - (viii). (i) and (ii) only
- (b) Suppose we define a vector `x <- c(3,0,-1,4,0,-2)`. What will `which(x==0)` give you?
- (i). A vector of the 0 elements of `x`.
 - (ii). A vectors of 0's.
 - (iii). A vector of `TRUE`'s and `FALSE`'s.
 - (iv). The vector of the indices of the 0 values.
-

Q3-6.

- (a) Define the matrix `A` as:

```
## Warning in matrix(c(1, 2, 3, 6, 1, 2, 4, 1, 1, 3, 3, 6, 4), nrow = 3): data
## length [13] is not a sub-multiple or multiple of the number of rows [3]
```

```
##      [,1] [,2] [,3] [,4] [,5]
## [1,]    1    6    4    3    4
## [2,]    2    1    1    3    1
## [3,]    3    2    1    6    2
```

What is the output of `apply(A, -1, 1, sd)`?

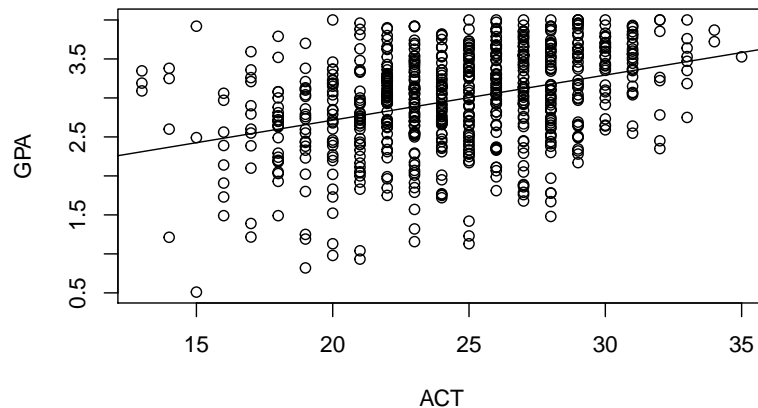
- (i). A vector of length 4 corresponding to the standard deviation of each column of `A`, excluding the first column.
 - (ii). A vector of length 3 corresponding to the standard deviation of each row of `A`, excluding the first column.
 - (iii). The standard deviation of all the values in `A`.
 - (iv). The standard deviation of the first row of `A`.
 - (v). An error since `A` doesn't have a -1 column.
- (b) Which of the following lines of code successfully constructs the matrix for part (a)? Comment on the strengths and weaknesses of the correct answers.
- (i). `cbind(c(1,2,3), c(6,1,2), c(4,1,1), c(3,3,6), c(4,1,2))`
 - (ii). `matrix(cbind(c(1,2,3), c(6,1,2), c(4,1,1), c(3,3,6), c(4,1,2)))`
 - (iii). `t(matrix(c(1,6,4,3,4,2,1,1,3,1,3,2,1,6,2),nrow = 3))`
 - (iv). `matrix(c(1,2,3,6,1,2,4,1,1,3,3,6,4),nrow = 3)`
-

Q4. Fitting a linear model by least squares

Q4-1.

The admissions officer at a large state university wants to assess how well academic success can be predicted based on information available at admission. She collects data on freshman GPA and highschool ACT exam scores for 705 students in an R dataframe called `gpa`. The plot below shows a line fitted to a scatterplot of the points in the dataset.

```
gpa_lm <- lm(GPA~ACT,data=gpa)
plot(GPA~ACT,data=gpa)
abline(coef(gpa_lm))
```



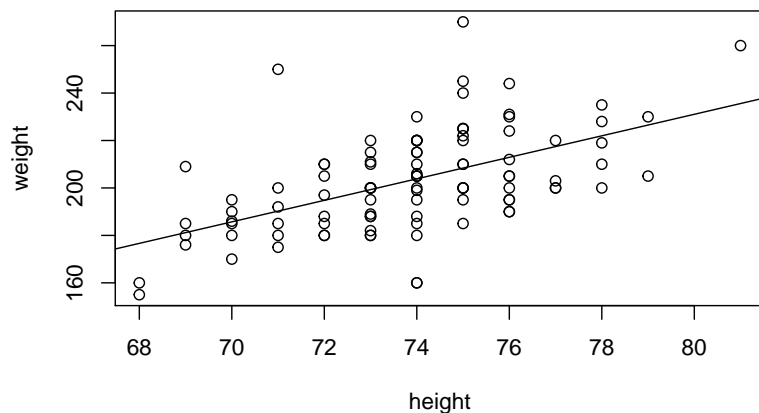
- (a) Explain in words the criterion that is used to obtain the fitted line in the plot above.
 - (b) Defining appropriate notation, write an equation for the fitted model in subscript form. At this point, you don't have to explain how the coefficients are calculated.
 - (c) Defining appropriate notation, write an equation for the fitted model in matrix form. You still don't have to explain how the coefficients are calculated.
 - (d) Now, explain using matrix notation how the model coefficients are calculated.
 - (e) Write an equation using subscript notation for the *fitted value* for the i th baseball player. Write a sentence to explain the interpretation of this fitted value.
-

Q4-2.

A statistician employed by a major league baseball team is asked to assess the range of typical weights for major league baseball players of a given height. She obtains data from http://wiki.stat.ucla.edu/socr/index.php/SOCR_Data_MLB_HeightsWeights and reads them into R as a dataframe including variables 'Height' (in inches) and 'Weight' (in pounds) for each of 1035 Major League Baseball players. She starts by analyzing just the first 100 players.

She fits a linear model and plots the data and the resulting fitted line using the following R code:

```
weight_lm <- lm(weight ~ height)
plot(height,weight)
abline(coef(weight_lm))
```



- (a) Write out the fitted linear model using subscript notation, including the following coefficients from `weight_lm`. This means you are asked to use actual numbers, rather than letters, for the model coefficients. Make sure to define any notation you introduce.

```
round(coef(weight_lm),3)
```

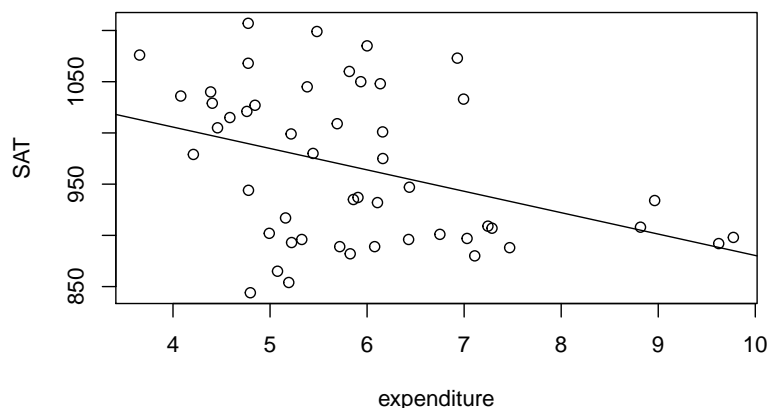
```
## (Intercept)      height
##    -131.652      4.534
```

- (b) Use matrix notation to explain how these coefficients were calculated.
- (c) The tenth observation corresponds to Adam Stern, an outfielder for the Baltimore Orioles. His recorded height is 71 inches. Write out the formula for the fitted value for this observation. You do not need to simplify your calculation.
- (d) Use matrix notation to write out an expression for the fitted values of the model. Make sure to define appropriate notation.

Q4-3.

The government wants to understand the relationship between expenditures on public education and test results. The dataset `SAT` contains the per-pupil annual expenditure (in thousands of dollars) and the average SAT score for each of the 50 states in 1994-95. The plot below shows a line fitted to a scatterplot of the points in the dataset.

```
sat_lm <- lm(SAT~expenditure,data=sat)
plot(SAT~expenditure,data=sat)
abline(coef(sat_lm))
```



- (a) Write out the regression model in subscript form (including an intercept term). Use letters, rather than actual numbers, for the model coefficients: you don't have the actual numbers at this point. Make sure to define any notation you introduce.

The table below shows the head of the dataset (the first 6 rows).

```
head(sat)
```

```
##      expenditure SAT
## Alabama      4.405 1029
## Alaska       8.963  934
## Arizona      4.778  944
## Arkansas     4.459 1005
## California   4.992  902
## Colorado     5.443  980
```

- (b) Write out the corresponding design matrix. You only need to put in actual numbers for the first 5 rows, use ... after that and specify the dimension of the matrix.
- (c) Explain how the model coefficients are evaluated. Name the method and give the appropriate formula.
- (d) Describe in one line what trend you observe from the plot (what would you interpret from this data). Is this what you would've expected? What could be a possible justification for the trend being observed?

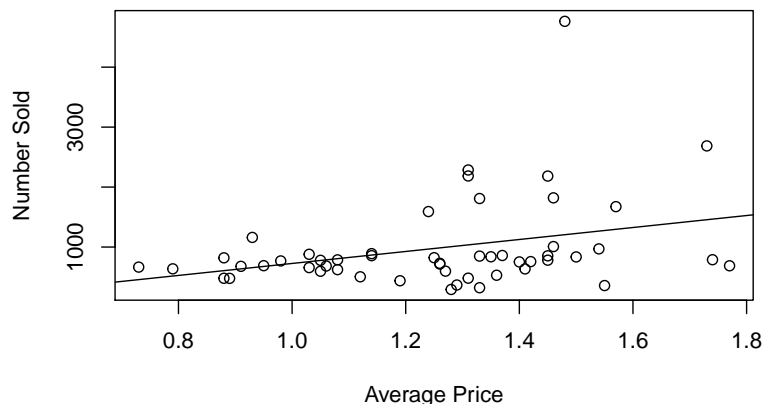
Q4-4.

A statistician employed by an avocado producer is asked to assess the relationship between avocado prices and sales volume for small Hass avocados. She obtains data from <https://www.kaggle.com/neuromusic/avocado-prices> and reads them into R as a dataframe. She keeps only the 2016 data for organic avocados sold in the Detroit area and plots the average price in dollars ('AveragePrice') against the number of small Hass avocados sold ('X4046'). This results in a dataset with 52 observations.

```
avocado <- read.csv("avocado.csv")
avocado_2016 <- subset(avocado, year == 2016 & type == 'organic' & region == 'Detroit')
```

She fits a linear model and plots the data and the resulting fitted line using the following R code:

```
price_lm <- lm(X4046 ~ AveragePrice, data = avocado_2016)
plot(avocado_2016$AveragePrice, avocado_2016$X4046,
     xlab = 'Average Price', ylab = 'Number Sold')
abline(coef(price_lm))
```



- (a) Write out the fitted linear model using subscript notation, including the following coefficients from `price_lm`. This means you are asked to use actual numbers, rather than letters, for the model coefficients. Make sure to define any notation you introduce.

```
round(coef(price_lm),3)
```

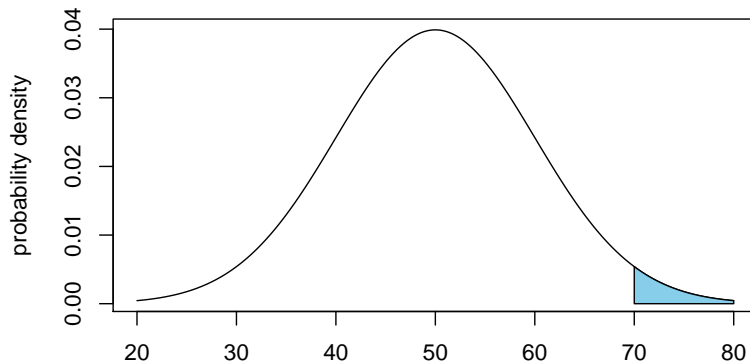
```
## (Intercept) AveragePrice
##      -272.829      999.062
```

- (b) Use matrix notation to explain how these coefficients were calculated.
- (c) Use matrix notation to write out an expression for the residual values of the model. Make sure to define appropriate notation.
- (d) From the scatter plot above, the statistician notices a potential outlier. This potential outlier corresponds to the the second week in January 2016. This week organic small Hass avocados sold for an average of \$1.48 with a total of 4,763 sold. Write a numeric expression for the residual of this observation—you are not expected to evaluate it.

Q5. Probability exercises

Q5-1.

The figure below shows the probability density function of a normal random variable X .



- (a) By looking at the probability density function, estimate the mean and standard deviation of X . Use these estimates for the subsequent parts of this question.
- (b) Write a probability statement about the random variable X that corresponding to the shaded area.
- (c) Write an integral corresponding to this shaded area.
- (d) Write R code to evaluate this integral numerically.

Q5-2.

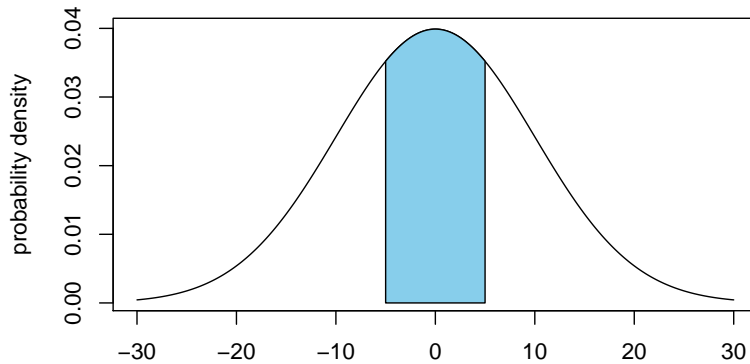
Let Y be a discrete random variable that takes values 0, 1, or 2 with probabilities 0.25, 0.5, and 0.25, respectively.

- (a) What is the expected value of Y ?

(b) What is the variance of Y ?

Q5-3.

The figure below shows the probability density function of a normal random variable X .



- (a) By looking at the probability density function, estimate the mean and standard deviation of X . Use these estimates for the subsequent parts of this question.
- (b) Write a probability statement about the random variable X that corresponding to the shaded area.
- (c) Write an integral corresponding to this shaded area.
- (d) Write R code to evaluate this integral numerically.

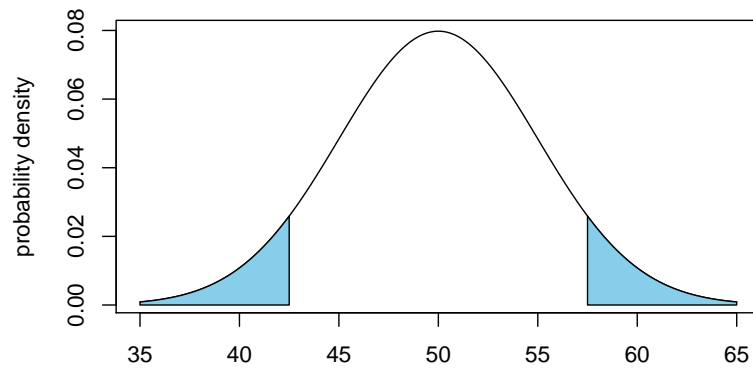
Q5-4.

The average midterm score of 20 students is 40 out of a maximum of 75. The professor realizes that she missed one student and upon including his score, the average went up by one point and became 41.

- (a) What is the midterm score of the 21st student?
- (b) Suppose the sample mean and variance are 41 and 36, and we model the midterm scores as being a draw from a normal distributed with these parameters. What is the chance that a student drawn at random gets over 59? Write your answer as a call to `pnorm()` and also give an approximate answer based on your knowledge of the normal distribution.

Q5-5.

The figure below shows the probability density function of a normal random variable X .



- (a) By looking at the probability density function, estimate the mean and standard deviation of X . Use these estimates for the subsequent parts of this question.
- (b) Write a probability statement about the random variable X corresponding to the shaded area.
- (c) Write an integral corresponding to this shaded area.
- (d) Write R code to evaluate this integral numerically.

License: This material is provided under an MIT license
