

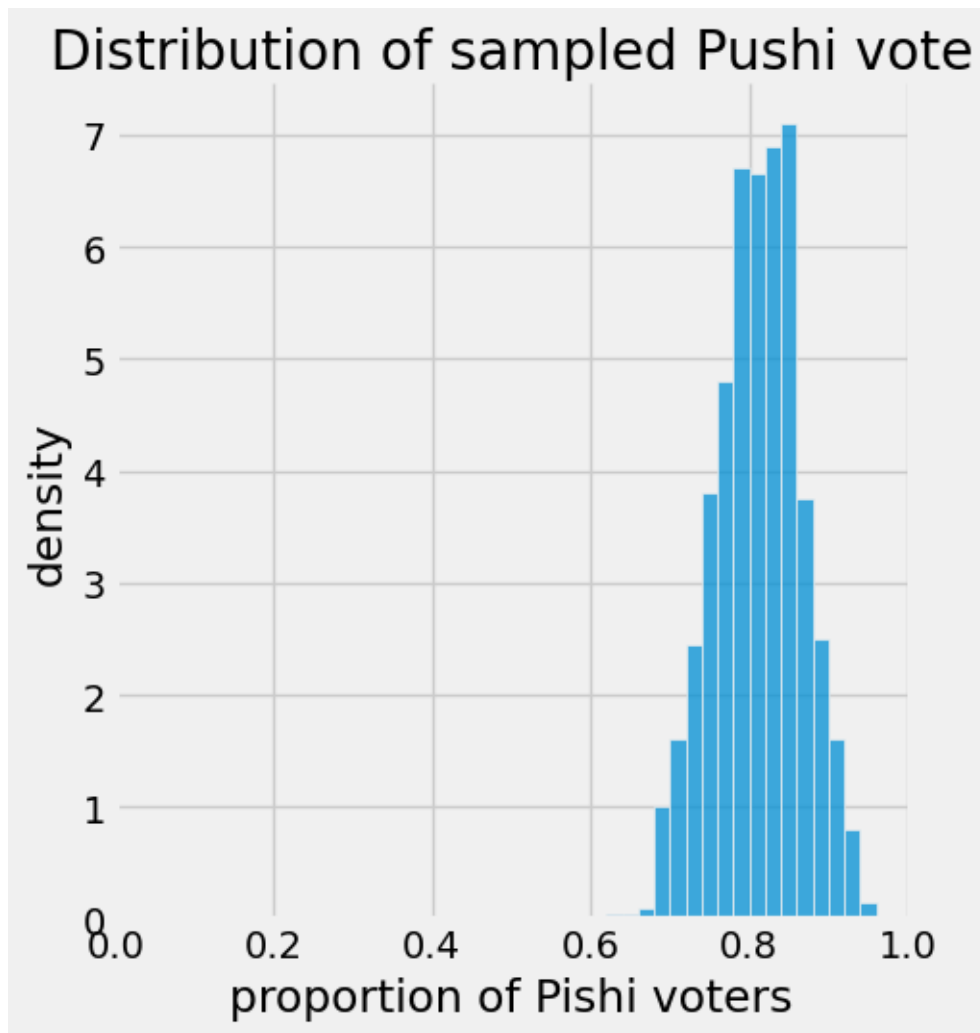
---

### 0.0.1 Question 2b

Create a plot using any `seaborn` and/or `matplotlib.pyplot` functions of your choice to visualize `samples`, which is the simulated distribution of Pishi votes using a sample of size 50. Include descriptive titles and labels. An example is included below. The total area under the plot must be normalized to 1. Your plot may not match exactly ours due to randomness of the data generating process in `np.random.multinomial`.

**Hint:** use `plt.xlim(left, right)` ([documentation](#)) to specify the left and right limits of the x-axis.

```
In [71]: sns.displot(samples, stat='density', binwidth=0.02)
         plt.xlim(0, 1)
         plt.title('Distribution of sampled Pushi vote')
         plt.xlabel('proportion of Pishi voters')
         plt.ylabel('density')
         plt.show()
```



---

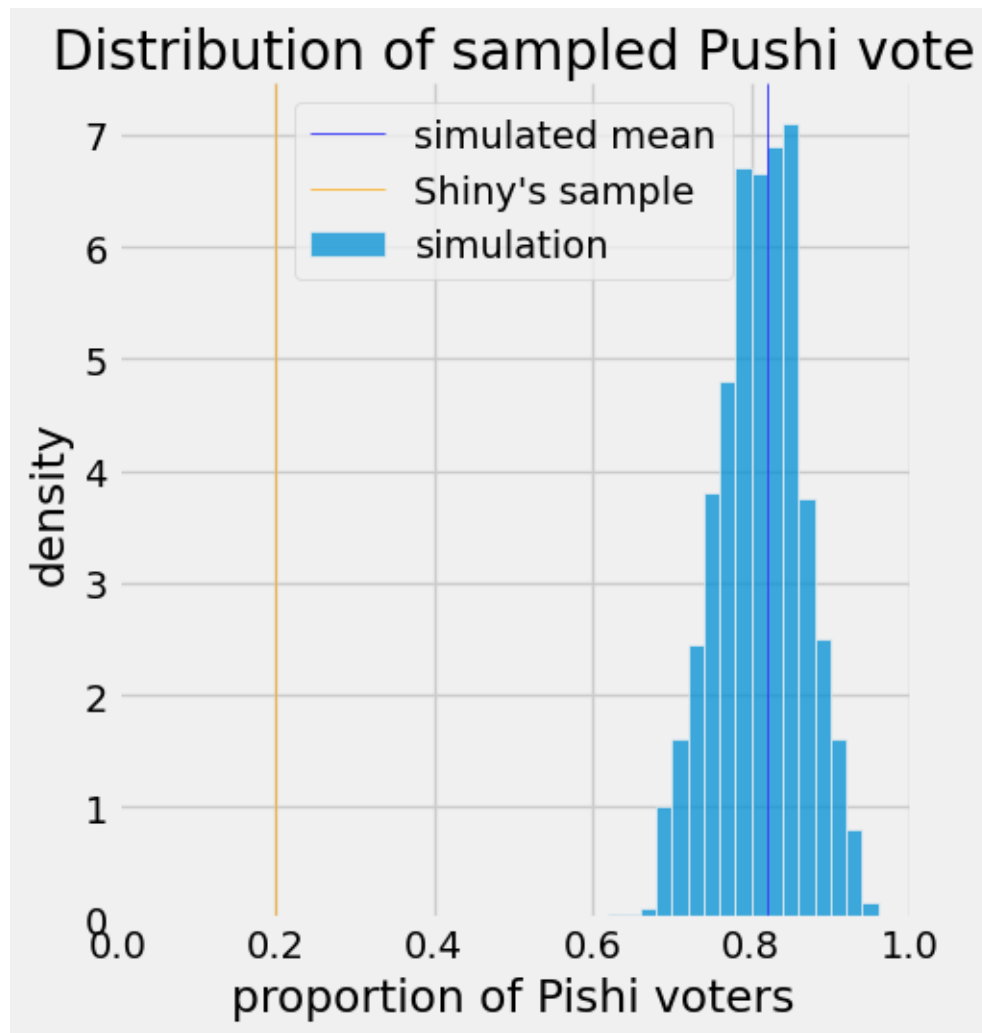
### 0.0.2 Question 2c

According to Shiny's 50-person sample, 20% of her discussion section reported that they would vote for Pishi in the end-of-semester contest.

In the cell below, create a plot using any `seaborn` and/or `matplotlib.pyplot` functions of your choice to visualize Shiny's sample statistic superimposed on the simulated sample distribution you plotted in the previous part. In other words, include - a vertical line that passes through 20%, - a vertical line that passes through the mean of the simulated sample distribution, and - the simulated sample distribution itself.

You should choose contrasting colors and include a descriptive title, labels, and a legend if needed. An example is included below.

```
In [77]: sns.displot(samples, stat='density', binwidth=0.02, label='simulation')
plt.xlim(0, 1)
plt.title('Distribution of sampled Pushi vote')
plt.xlabel('proportion of Pishi voters')
plt.ylabel('density')
plt.axvline(x=np.mean(samples), color='blue', label='simulated mean', lw=0.5)
plt.axvline(x=0.2, color='orange', label="Shiny's sample", lw=0.5)
plt.legend()
plt.show()
```



---

### 0.0.3 Question 2d

Based on your analysis above, could Shiny's result have arisen due to chance alone? If not, what could be a potential source of bias?

No, even in the most severe outliers, Shiny's result could have not arisen due to chance alone. Since Shiny was using a convenience sample (her discussion session), this could lead to a potential bias in people voting for Mimi (by far the majority in her sample) just because the owner is Shiny (the discussion GSI). Because of the chance that this could happen, this could be the source of bias, and would most likely make Shiny's sample not a good sample to take from.



## Homework #5B

**Total Points: 24**

### Submission Instructions

You must submit this assignment to Gradescope by the on-time deadline, **Thursday, October 5th at 11:59 PM Pacific**. Please read the syllabus for **the grace period policy**. No late submissions beyond the grace period will be accepted. While course staff is happy to help you if you encounter difficulties with submission, we may not be able to respond to last-minute requests for assistance (TAs need to sleep, after all!). **We strongly encourage you to plan to submit your work to Gradescope several hours before the stated deadline.** This way, you will have ample time to reach out to staff for submission support.

There are two parts to this assignment listed on Gradescope:

- **Homework 05 Coding:** Submit your Jupyter notebook zip file for Homework 5A, which can be generated and downloaded from DataHub by using the `grader.export()` cell provided.
- **Homework 05 Written:** Submit a single PDF to Gradescope that contains both (1) your answers to all manually graded questions from the Homework 5A Jupyter Notebook and (2) your answers to all questions in this Homework 5B document.

To receive credit on this assignment, **you must submit both your coding and written portions to their respective Gradescope portals**. Your written submission (a single PDF) can be generated as follows:

1. Access your answers to manually graded Homework 5A questions in one of three ways:
  - *Automatically create PDF (recommended):* We have provided a cell to generate your written response in the Homework 5A notebook for you. Run the cell and click to download the generated PDF. This function will extract your response to the manually-graded questions and put them on separate pages. This process may fail if your answer is not properly formatted; if this is the case, check out common errors and solutions described on Ed or follow either of the two ways described below.

- *Manually download PDF*: If there are issues with automatically generating the PDF, on DataHub, you can try downloading the PDF by clicking on **File->SaveandExportNotebookAs...->PDF**. If you choose to go this route, you must take special care to ensure all appropriate pages are chosen for each question on Gradescope.
  - *Take screenshots*: If that doesn't work either, you can take screenshots of your answers (and your code if present) to manually-graded questions and include them as images in a PDF. The manually-graded questions are listed at the top of the Homework 5A notebook.
2. Answer the below Homework 5B written questions in one of many ways:
    - You can type your answers. We recommend LaTeX, the math typesetting language. Overleaf is a great tool to type in LaTeX.
    - Download this PDF, print it out, and write directly on these pages. If you have a tablet, you may save this PDF and write directly on it.
    - Write your answers on a blank sheet of physical or digital paper.
    - Note: If you write your answers on physical paper, use a scanning application (e.g., CamScanner, Apple Notes) to generate a PDF.
  3. Combine these two sets of answers together into one PDF document and submit it to the appropriate Gradescope written portal. You can use PDF merging tools, e.g., Adobe Reader, Smallpdf (<https://smallpdf.com/merge-pdf>) or Apple Preview (<https://support.apple.com/en-us/HT202945>).
  4. **Important**: When submitting on Gradescope, you **must tag pages to each question correctly** (it prompts you to do this after submitting your work). This significantly streamlines the grading process for our readers. Failure to do this may result in a score of 0 for untagged questions.

*You are responsible for ensuring your submission follows our requirements. We will not be granting regrade requests nor extensions to submissions that don't follow instructions.* If you encounter any difficulties with submission, please don't hesitate to reach out to staff prior to the deadline.

## Collaborators

Data science is a collaborative activity. While you may talk with others about the homework, we ask that you write your solutions individually. If you do discuss the assignments with others, please include their names at the top of your submission.



## Properties of Linear Regression Residuals

1. (10 points) In the lecture, we spent a great deal of time talking about Simple Linear Regression (SLR), which you also saw in Data 8. To briefly summarize, the simple linear regression model assumes that given a single observation  $x$ , our predicted response for this observation is  $\hat{y} = \theta_0 + \theta_1 x$ .

In Lecture 10, we saw that the  $\theta_0 = \hat{\theta}_0$  and  $\theta_1 = \hat{\theta}_1$  that minimize the average  $L_2$  loss (or Mean Squared Error - MSE) for the simple linear regression model are:

$$\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}$$

$$\hat{\theta}_1 = r \frac{\sigma_y}{\sigma_x}$$

Or, rearranging terms, our predictions  $\hat{y}$  are:

$$\hat{y} = \bar{y} + r \sigma_y \frac{x - \bar{x}}{\sigma_x}$$

- (a) (3 points) As we saw in the lecture, a residual  $e_i$ , for data point  $i \in \{1, \dots, n\}$ , is defined to be the difference between a true response  $y_i$  and predicted response  $\hat{y}_i$ . Specifically,  $e_i = y_i - \hat{y}_i$ . Note that there are  $n$  data points, and each data point is denoted by  $(x_i, y_i)$ .

Prove, using the equation for  $\hat{y}$  above, that  $\sum_{i=1}^n e_i = 0$ .

Handwritten work for part (a):

$$\sum e_i \leftarrow \begin{aligned} & \sum (y_i - \hat{y}_i) \\ & = \sum y_i - \sum \left( \bar{y} + r \sigma_y \frac{x_i - \bar{x}}{\sigma_x} \right) \\ & = \sum y_i - n\bar{y} - r \sigma_y \frac{\sum x_i - n\bar{x}}{\sigma_x} \end{aligned}$$

Since  $\sum x_i = n\bar{x}$ , the last term is zero.

$$\sum e_i = \sum y_i - n\bar{y} = 0$$

- (b) (2 points) Prove that  $\bar{\hat{y}} = \bar{y}$ . You may use your result from part (a).

Handwritten work for part (b):

$$\bar{\hat{y}} = \frac{1}{n} \sum (\bar{y} + r \sigma_y \frac{x_i - \bar{x}}{\sigma_x})$$

$$\bar{\hat{y}} = \bar{y} + r \sigma_y \frac{\sum x_i - n\bar{x}}{n\sigma_x}$$

Since  $\sum x_i = n\bar{x}$ , the last term is zero.

$$\bar{\hat{y}} = \bar{y}$$

- (c) (2 points) Show that  $(\bar{x}, \bar{y})$  is on the simple linear regression line.

Handwritten work for part (c):

$$\hat{y} = \bar{y} + r \sigma_y \frac{x_i - \bar{x}}{\sigma_x}$$

$$= \bar{y} + 0$$

$$\hat{y} = \bar{y} \text{ only when } x_i = \bar{x},$$

$$\therefore (\bar{x}, \bar{y}) \text{ is simple lin reg.}$$

- (d) (3 points) Show that the residuals are uncorrelated with the predictor variable, that is

$$\frac{1}{n} \sum_{i=1}^n \left( \frac{e_i - \bar{e}}{\sigma_e} \right) \left( \frac{x_i - \bar{x}}{\sigma_x} \right) = 0,$$

where  $\bar{e} = \frac{1}{n} \sum_{i=1}^n e_i$ ,  $\sigma_e^2 = \frac{1}{n} \sum_{i=1}^n (e_i - \bar{e})^2$ , and  $\sigma_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ . You may assume that  $\sigma_e$ ,  $\sigma_x$ , and at least one residual are not exactly zero. Use the properties of estimating equations derived in the lecture.

$\bar{e} = \text{mean of residuals} = 0$

$$\therefore \frac{1}{n} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sigma_x} \right) \left( \frac{e_i}{\sigma_e} \right) = 0$$

$$\frac{1}{n} \sum \frac{x_i e_i}{\sigma_x \sigma_e} - \sum \frac{\bar{x} e_i}{\sigma_e \sigma_x} = 0$$

$$\frac{1}{n} \sum \frac{x_i e_i}{\sigma_x \sigma_e} - \bar{x} \sum \frac{e_i}{\sigma_e \sigma_x} = 0$$

$\downarrow$   
 $\bar{e} = 0$

$$\frac{1}{n} \sum \frac{x_i e_i}{\sigma_x \sigma_e} = 0$$

$$\left( \frac{1}{\sigma_x \sigma_e} \sum x_i e_i = 0 \right) n \sigma_x \sigma_e$$

$$\sum x_i e_i = 0$$

$$\text{Covariance} = 0$$

## Properties of a Linear Model With No Constant Term

2. (4 points) Suppose that we don't include an intercept term in our model. That is, our model is now

$$\hat{y} = \theta x,$$

where  $\theta$  is the single parameter for our model that we need to optimize. (In this equation,  $x$  is a scalar, corresponding to a single observation.)

As usual, we are looking to find the value  $\hat{\theta}$  that minimizes the average  $L_2$  loss (MSE) across our observed data  $\{(x_i, y_i)\}$ , for  $i \in \{1, \dots, n\}$ :

$$R(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \theta x_i)^2$$

The estimating equations derived in the lecture no longer hold. In this problem, we'll derive a solution to this simpler model. We'll see that the least squares estimate of the slope in this model differs from the simple linear regression model, and we'll also explore whether or not our properties from the previous problem still hold.

Use calculus to find the minimizing  $\hat{\theta}$ .

That is, you may prove that:

$$\hat{\theta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

Hint: You may start by following the format of SLR in lecture 10 and replace the SLR model with the model defined above.

$$R'(\theta) = -\frac{1}{n} \sum (2y_i - 2\theta x_i)(x_i)$$

$$\left. \begin{array}{l} (y_i^2 - 2y_i \theta x_i + \theta^2 x_i^2) \\ (2y_i - 2\theta x_i)(x_i) \end{array} \right\}$$

$$\sum_{i=1}^n \theta x_i^2 = \sum_{i=1}^n y_i x_i$$

$$\hat{\theta} \sum x_i^2 = \sum y_i x_i$$

## MSE “Minimizer”

3. (10 points) Recall from calculus that given some function  $g(x)$ , the  $x$  you get from solving  $\frac{dg(x)}{dx} = 0$  is called a *critical point* of  $g$  – this means it could be a minimizer or a maximizer for  $g$ . In this question, we will explore some basic properties and build some intuition on why, for certain loss functions such as squared  $L_2$  loss, the critical point of the empirical risk function (defined as an average loss on the observed data) will always be the minimizer.

Given some linear model  $f(x) = \theta x$  for some real scalar  $\theta$ , we can write the empirical risk of the model  $f$  given the observed data  $\{x_i, y_i\}$ , for  $i \in \{1, \dots, n\}$  as the average  $L_2$  loss (MSE):

$$\frac{1}{n} \sum_{i=1}^n (y_i - \theta x_i)^2 = \sum_{i=1}^n \frac{1}{n} (y_i - \theta x_i)^2$$

- (a) (3 points) Let's investigate one of the  $n$  functions in the summation in the MSE. Define  $g_i(\theta) = \frac{1}{n}(y_i - \theta x_i)^2$  for  $i \in \{1, \dots, n\}$ . In this case, note that the MSE can be written as  $\sum_{i=1}^n g_i(\theta)$ .

Recall from calculus that we can use the 2nd derivative of a function to describe its curvature about a certain point (if it is facing concave up, down, or possibly a point of inflection). You can take the following as a fact: A function is convex if and only if the function's 2nd derivative is non-negative on its domain. Based on this property, verify that  $g_i(\theta)$  is a **convex function**.

$$\begin{aligned} g(\theta) &= \frac{1}{n} (y_i - \theta x_i)^2 \\ &= \frac{1}{n} (y_i^2 - 2\theta x_i y_i + \theta^2 x_i^2) \\ &= -\frac{1}{n} (2y_i x_i - \theta^2 x_i^2) \end{aligned} \quad \begin{aligned} g''(\theta) &= \frac{1}{n} (2x_i^2) \\ g_i(\theta) &\text{ is a convex function} \end{aligned}$$

- (b) (2 points) Briefly explain intuitively in words why given a convex function  $g(\theta)$ , the critical point we get by solving  $\frac{dg(\theta)}{d\theta} = 0$  minimizes  $g$ . You can assume that  $\frac{dg(\theta)}{d\theta}$  is a function of  $\theta$  (and not a constant).

If  $g$  is a convex function it goes from  $-$  to  $+$  in  $g'(\theta)$ , setting it to 0 means you catch at the lowest point ~~the~~  $g'(\theta) = 0$

- (c) (3 points) Now that we have shown that each term in the summation of the MSE is a convex function, one might wonder if the entire summation is convex, given

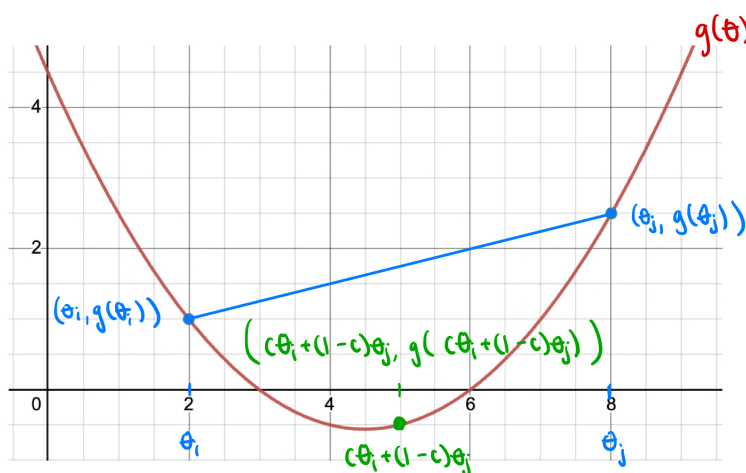
that it is a sum of convex functions.

Let's look at the formal definition of a **convex function**. Algebraically speaking, a function  $g(\theta)$  is convex if for any two points  $(\theta_i, g(\theta_i))$  and  $(\theta_j, g(\theta_j))$  on the function,

$$g(c \times \theta_i + (1 - c) \times \theta_j) \leq c \times g(\theta_i) + (1 - c) \times g(\theta_j)$$

for any real constant  $0 \leq c \leq 1$ .

The function  $g$  evaluated on any point between  $\theta_i$  and  $\theta_j$  will always lie at or below the secant line connecting  $g(\theta_i)$  and  $g(\theta_j)$



See a graph in this Wikipedia article [https://en.wikipedia.org/wiki/Convex\\_function](https://en.wikipedia.org/wiki/Convex_function).

Intuitively, the above definition says that, given the plot of a convex function  $g(\theta)$ , if you connect 2 randomly chosen points on the function, the line segment will always lie on or above  $g(\theta)$  (try this with the graph of  $g(\theta) = \theta^2$ ).

- i. (2 points) Using the definition above, show that if  $g(\theta)$  and  $h(\theta)$  are both convex functions, their sum  $g(\theta) + h(\theta)$  will also be a convex function.

$$\begin{aligned} g(c\theta_i + (1-c)\theta_j) &\leq c g(\theta_i) + (1-c) g(\theta_j) \\ + h(c\theta_i + (1-c)\theta_j) &\leq c h(\theta_i) + (1-c) h(\theta_j) \\ &= (g+h)(c\theta_i + (1-c)\theta_j) \leq c(g+h)(\theta_i) + (1-c)(g+h)(\theta_j) \end{aligned}$$

Since the inequality does not change after adding, has to be convex

- ii. (1 point) Based on what you have shown in the previous part, explain intuitively why a (finite) sum of  $n$  convex functions is still a convex function when  $n > 2$ .

if we proved that  $h + g$  was convex, then treat that as a new function and add a 3rd, or a 4th, or infinite.  $((g+h)+f)+\dots$

(d) (2 points) Remember from part (a) that the MSE can be written as:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \theta x_i)^2 = \sum_{i=1}^n \frac{1}{n} (y_i - \theta x_i)^2 = \sum_{i=1}^n g_i(\theta)$$

Explain why solving for the critical point of the MSE by taking the gradient with respect to the parameter  $\theta$  and setting that expression to 0, is guaranteed that the solution we find will minimize the MSE.

$(y_i - \theta x_i)^2$  is convex.

can use in MSE due to the convexity

in a) for a set  $n \{1, 2, \dots, n\}$ ; due to these factors,   
 Closing note: In this question, we have discussed only the simple linear model with no constant term—a single-variable function. However, the above properties extend more generally to all multivariable linear regression models; this proof is beyond the scope of this course and is left to a future you.   
 guarantees minimizes

**Congratulations! You have finished Homework 5B!**