

## Discussion #5

Name:

## Dimensionality Reduction

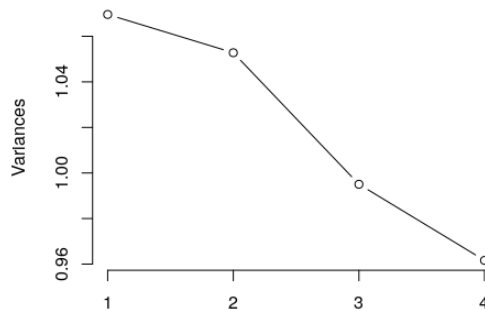
- Principal Component Analysis (PCA) is one of the most popular dimensionality reduction techniques because it is relatively easy to compute and its output is interpretable. To get a better understanding of what PCA is doing to a dataset, let's imagine applying it to points contained within this surfboard. The origin is in the center of the board, and each point within the board has three attributes: how far (in inches) along the board's length, width, and thickness the point is from the center. These three dimensions determine the spread of the data.



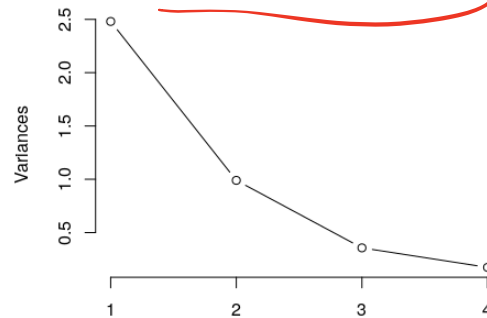
- If we were to apply PCA to the surfboard, what would the first three principal components (PCs) represent? Feel free to draw and label these dimensions on the image of the surfboard.
  - Which of the three PCs should be used to create a 2D representation of the surfboard? How come? Make a sketch of the 2D projection below. *length, width*
- Compare the scree plots produced by performing PCA on dataset A and on dataset B. For which of the datasets would PCA provide a scatter plot that describes the variability of the data without leaving out much information? Note that the columns of both datasets were centered to have means of 0 and scaled to have a variance of 1.

$$\text{singular values of } X = \sqrt{\text{eigenvalues of } X^T X}$$

Scree Plot of Dataset A



Screeplot of Dataset B



most of the information is conveyed by PC1, PC2

## Midterm Review

### 1. Probability and Sampling

3. A small town has 5 houses with the following people living in each house:



Suppose we take a **cluster sample** of 2 houses (without replacement), what is the chance that:

(a) Kim and Lars are in the sample

- ☐ 0    ☐ 1/20    ☐ 1/10    ☐ 1/6    ☐ 1/5    ☐ 2/5    ☐ 1

You may show your work in the following box for partial credit:

$$p = \frac{\text{\# samples w/ house (5)}}{\text{total \# of samples}} = \frac{\binom{4}{1}}{\binom{5}{2}} = \frac{4}{10} = \frac{2}{5}$$

(b) Kim, Abe, and Ben are in the sample

$$\frac{5!}{3!2!} = \frac{5 \cdot 4 \cdot 3!}{3! \cdot 2!} = 10$$

$\binom{n}{k}$ :  
# of ways to select k objects from a pool of n

☐ 0   
 ☐ 1/20   
 ☐ 1/10   
 ☐ 1/6   
 ☐ 1/5   
 ☐ 2/5   
 ☐ 1

You may show your work in the following box for partial credit:

$$\begin{aligned}
 p(\text{house } \textcircled{1} \text{ and house } \textcircled{5}) &= \frac{1}{\binom{5}{2}} = \frac{1}{10} \\
 \textcircled{5} \frac{1}{5} \cdot \textcircled{1} \frac{1}{4} + \textcircled{1} \frac{1}{5} \cdot \textcircled{5} \frac{1}{4} &= \frac{1}{10}
 \end{aligned}$$

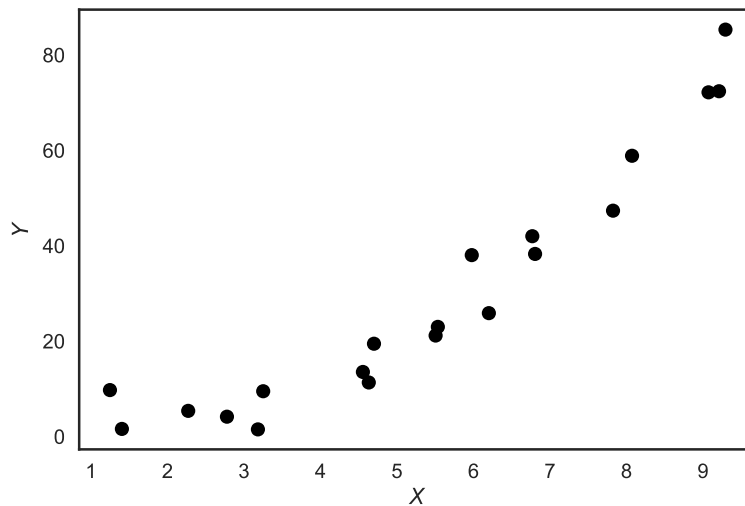
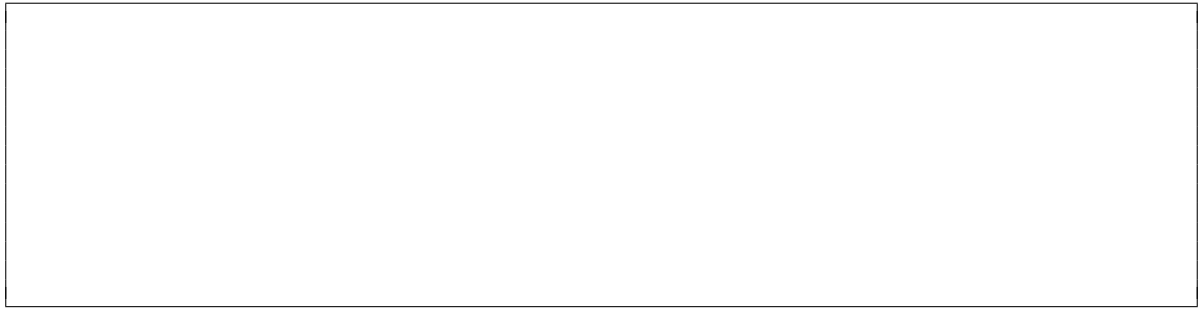
(c) Kim and Dan are in the sample - <sup>1/10</sup> **Select all that apply**

- ☐ The same as the chance Kim and Lars are in the sample <sup>2/5</sup>  
☒ The same as the chance Kim, Abe, and Ben are in the sample <sup>1/10</sup>  
☐ Neither of the above

## 2. Transformations and Smoothing

4. Which of the following are reasonable motivations for applying a power transformation? **Select all that apply:**

- ☒ To help visualize highly skewed distributions  
☒ Bring data distribution closer to random sampling  
☒ To help straighten relationships between pairs of variables.  
☒ Reduce the dimension of data  
☒ Remove missing values



$$y \approx x^2$$

5. Which of the following transformations could help make linear the relationship shown in the plot below? **Select all that apply:**

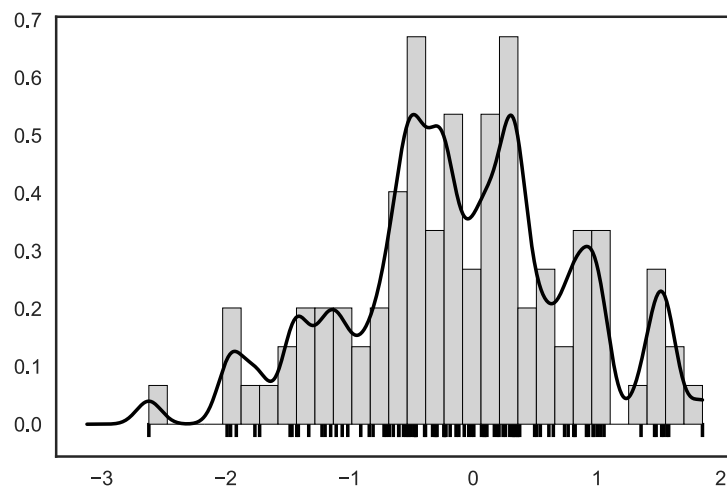
☐  $\log(y)$    ☐  $x^2$    ☐  $\sqrt{y}$    ☐  $\log(x)$    ☐  $y^2$    ☐ None of the above



6. The above plot contains a histogram, rug plot, and Gaussian kernel density estimator. The Gaussian kernel is defined by:

$$K_{\alpha}(x, z) = \frac{1}{\sqrt{2\pi\alpha^2}} \exp\left(-\frac{(x - z)^2}{2\alpha^2}\right)$$

Judging from the shape of separate standing peaks, which of the following is the most likely value for the kernel parameter  $\alpha$ .



☐  $\alpha = 0$    ☒  $\alpha = 0.1$    ☐  $\alpha = 10$    ☐  $\alpha = 100$