Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

# Discussion 12

Exam Review

April 26, 2018

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

① Probability & Sampling

② EDA & Visualization

③ Prediction

④ Optimization

⑤ Inference

⑥ Big Data

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

A political scientist is interested in answering a question about a country composed of three states with exactly 10000, 20000, and 30000 voting adults. To answer this question, a political survey is administered by randomly sampling 25, 50, and 75 voting adults from each town in each state, respectively.

**Which sampling plan was used in the survey?**

(a) cluster sampling

(b) stratified sampling

(c) quota sampling

(d) census

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

A political scientist is interested in answering a question about a country composed of three states with exactly 10000, 20000, and 30000 voting adults. To answer this question, a political survey is administered by randomly sampling 25, 50, and 75 voting adults from each town in each state, respectively. **Which sampling plan was used in the survey?**

(b) stratified sampling

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Suppose Sam visits your store to buy some items. He buys toothpaste for \$2.00 with probability 0.5. He buys a toothbrush for \$1.00 with probability 0.1. Let the random variable $X$ be the total amount Sam spends. Find $\mathbb{E}[X]$.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Suppose Sam visits your store to buy some items. He buys toothpaste for \$2.00 with probability 0.5. He buys a toothbrush for \$1.00 with probability 0.1. Let the random variable $X$ be the total amount Sam spends. Find $\mathbb{E}[X]$.

Let $X_{\text{toothpaste}}$ be the amount Sam spends on toothpaste, and $X_{\text{toothbrush}}$ be the amount Sam spends on a toothbrush. From the linearity of expectation, we have:

$$\mathbf{E}[X] = \mathbf{E}[X_{\text{toothpaste}} + X_{\text{toothbrush}}] = \mathbf{E}[X_{\text{toothpaste}}] + \mathbf{E}[X_{\text{toothbrush}}]$$

We know that $\mathbf{E}[X_{\text{toothpaste}}] = (0.5)(0) + (0.5)(2) = 1$, and $\mathbf{E}[X_{\text{toothbrush}}] = (0.9)(0) + (0.1)(1) = 0.1$. Thus, $\mathbf{E}[X] = 1.1$.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Suppose we have a coin that lands heads 80% of the time. Let the random variable $X$ be the *proportion* of times the coin lands tails out of 100 flips. What is $\text{Var}[X]$?

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Suppose we have a coin that lands heads 80% of the time. Let the random variable $X$ be the *proportion* of times the coin lands tails out of 100 flips. What is $\text{Var}[X]$?

Let $X_i$ be the outcome of the $i^{\text{th}}$ spin. If the $i^{\text{th}}$ spin lands heads than we say $X_i = 1$ and otherwise $X_i = 0$. Then the *proportion of times $X_i$ lands heads* is given by:

$$Y = \frac{1}{100} \sum_{i=1}^{n} X_i$$

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

We can compute the variance of $Y$ using the following identities:

$$\mathbf{Var}\left[Y\right] = \mathbf{Var}\left[\frac{1}{100}\sum_{i=1}^{n}X_i\right] \tag{1}$$

$$= \frac{1}{100^2}\mathbf{Var}\left[\sum_{i=1}^{n}X_i\right]$$

(Squared variance of constant multiple.)

$$= \frac{1}{100^2}\sum_{i=1}^{n}\mathbf{Var}\left[X_i\right]$$

(Ind. Variables implies linearity of var.)

$$= \frac{1}{100^2}\sum_{i=1}^{n}p(1-p) = \frac{p(1-p)}{100}$$

$$= \frac{.8(1-.8)}{100} = \frac{.16}{100} = .0016$$

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

For each of the following scenarios, determine which plot type
is *most* appropriate to reveal the distribution of and/or the
relationships between the following variable(s). For each
scenario, select only one plot type. Some plot types may be
used multiple times.

A. histogram

B. pie chart

C. bar plot

D. line plot

E. side-by-side boxplots

F. scatter plot

G. stacked bar plot

H. overlaid line plots

Sale price and number of bedrooms for houses sold in Berkeley in 2010.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Sale price and number of bedrooms for houses sold in Berkeley in 2010.

**E. Side-by-side Boxplots.**

We might imagine using a scatter plot since we are plotting the relationship between two numeric quantities. However because the number of bedrooms is an integer and most houses will only have a small number, we are likely to encounter *over-plotting* in the scatter plot. Therefore side-by-side boxplots are likely to be most informative.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Sale price and date of sale for houses sold in Berkeley between 1995 and 2015.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Sale price and date of sale for houses sold in Berkeley between 1995 and 2015.
**F. Scatter Plot.**
Here we are plotting two numeric quantities with sufficient spread on each axis.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Infant birth weight (grams) for babies born at Alta Bates
hospital in 2016.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Infant birth weight (grams) for babies born at Alta Bates
hospital in 2016.
**A. Histogram.**
Here we are plotting the distribution of a likely large number
of observations and therefore a histogram would be most
appropriate.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Mother's education-level (highest degree held) for students admitted to UC Berkeley in 2016.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Mother's education-level (highest degree held) for students admitted to UC Berkeley in 2016.
**C. Bar Plot.** Here we want to visualize counts of a categorical variable.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

SAT score and HS GPA of students admitted to UC Berkeley in 2016.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

SAT score and HS GPA of students admitted to UC Berkeley in 2016.
**F. Scatter Plot.** Here we are visualizing the relationship between two continuous quantities.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

The percentage of female student admitted to UC Berkeley each year from 1950 to 2000.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

The percentage of female student admitted to UC Berkeley each year from 1950 to 2000.
**D. Line plot.**
This allows us to see the trends over time.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

SAT score for males and females of students admitted to UCB
from 1950 to 2000

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

SAT score for males and females of students admitted to UCB
from 1950 to 2000
**E. side-by-side boxplots.**
This allows us to see the distributions of SAT scores per gender
and year.

Discussion 12

Exam Review

Probability & Sampling

EDA & Visualization

Prediction

Optimization

Inference

Big Data

When developing a model for a donkey's weight, we consider the following box plots of weight by age category.



This plot suggests:
- (a) Age is not needed in the model
- (b) Some of the age categories can be combined
- (c) Age could be treated as a numeric variable
- (d) None of the above

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

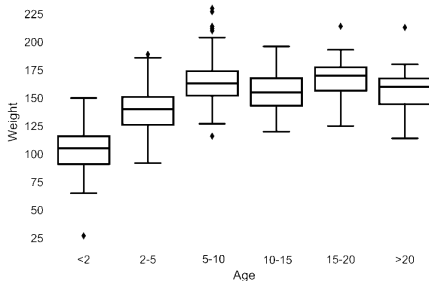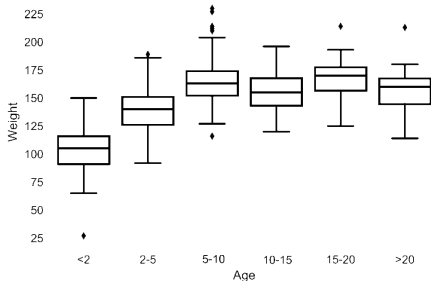When developing a model for a donkey's weight, we consider the following box plots of weight by age category.



This plot suggests:

(b) Some of the age categories can be combined

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Fix the following buggy Python implementation of gradient descent:

```python
def grad_descent(X, Y, theta0, grad_function,
    max_iter = 1000):
    """X: A 2D array, the feature matrix.
    Y: A 1D array, the response vector.
    theta0: A 1D array, the initial parameter
        vector.
    grad_function: Maps a parameter vector, a
        feature matrix, and a response vector to
        the gradient of some loss function at the
        given parameter value. The return value
        is a 1D array."""
    theta = theta0
    for t in range(1, max_iter+1):
        grad = grad_function(theta, X, Y)
        theta = theta0 + t * grad
    return grad
```

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

The last two lines need to change:

```
1  def grad_descent(X, Y, theta0, grad_function,
     max_iter = 1000):
2      """X: A 2D array, the feature matrix.
3      Y: A 1D array, the response vector.
4      theta0: A 1D array, the initial parameter
         vector.
5      grad_function: Maps a parameter vector, a
         feature matrix, and a response vector to
         the gradient of some loss function at the
         given parameter value. The return value
         is a 1D array."""
6      theta = theta0
7      for t in range(1, max_iter+1):
8          grad = grad_function(theta, X, Y)
9          theta = theta - (1/t) * grad
10     return theta
```

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Suppose you are given a dataset $\{(x_i, y_i)\}_{i=1}^{n}$ where $x_i \in \mathbb{R}$ is a one dimensional feature and $y_i \in \mathbb{R}$ is a real-valued response. You use $f_\theta$ to model the data where $\theta$ is the model parameter. You choose to use the following regularized loss:

$$L(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - f_\theta(x_i) \right)^2 + \lambda \theta^2$$

You choose to use the following regularized loss:

$$L(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - f_\theta(x_i) \right)^2 + \lambda \theta^2$$

This regularized loss is best described as:

(a) Average absolute loss with $L^2$ regularization.

(b) Average squared loss with $L^1$ regularization.

(c) Average squared loss with $L^2$ regularization.

(d) Average Huber loss with $\lambda$ regularization.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

You choose to use the following regularized loss:

$$L(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left(y_i - f_\theta(x_i)\right)^2 + \lambda\theta^2$$

This regularized loss is best described as:

(c) Average squared loss with $L^2$ regularization.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Suppose you choose the model $f_\theta(x_i) = \theta x_i^3$. Using the above
objective derive the loss minimizing estimate for $\theta$.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

**Step 1:** Take the derivative of the loss function.

$$\frac{\partial}{\partial \theta} L(\theta) = \frac{1}{n} \sum_{i=1}^{n} \frac{\partial}{\partial \theta} \left( y_i - \theta x_i^3 \right)^2 + \frac{\partial}{\partial \theta} \lambda \theta^2 \qquad (2)$$

$$= -\frac{2}{n} \sum_{i=1}^{n} \left( y_i - \theta x_i^3 \right) x_i^3 + 2\lambda\theta \qquad (3)$$

**Step 2:** Set derivative equal to zero and solve for $\theta$.

$$0 = -\frac{2}{n} \sum_{i=1}^{n} \left( \gamma_i - \theta x_i^3 \right) x_i^3 + 2\lambda\theta \qquad (4)$$

$$\theta = \frac{1}{n\lambda} \sum_{i=1}^{n} \left( \gamma_i - \theta x_i^3 \right) x_i^3 \qquad (5)$$

$$\theta = \frac{1}{n\lambda} \sum_{i=1}^{n} \gamma_i x_i^3 - \theta \frac{1}{n\lambda} \sum_{i=1}^{n} x_i^6 \qquad (6)$$

$$\theta \left( 1 + \frac{1}{n\lambda} \sum_{i=1}^{n} x_i^6 \right) = \frac{1}{n\lambda} \sum_{i=1}^{n} \gamma_i x_i^3 \qquad (7)$$

Discussion 12

Exam Review

Probability & Sampling

EDA & Visualization

Prediction

Optimization

Inference

Big Data

$$\theta \left( 1 + \frac{1}{n\lambda} \sum_{i=1}^{n} x_i^6 \right) = \frac{1}{n\lambda} \sum_{i=1}^{n} \gamma_i x_i^3 \tag{8}$$

Thus we obtain the final answer:

$$\boxed{\hat{\theta} = \frac{\frac{1}{n} \sum_{i=1}^{n} \gamma_i x_i^3}{\left( \lambda + \frac{1}{n} \sum_{i=1}^{n} x_i^6 \right)}} \tag{9}$$

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

**True or False.**
Suppose we have 100 samples drawn independently from a population. If we construct a 95% confidence interval for each sample, we expect 95 of them to include the **sample** mean.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

**True or False.**
Suppose we have 100 samples drawn independently from a
population. If we construct a 95% confidence interval for each
sample, we expect 95 of them to include the **sample** mean.
**False.** All of them should include the sample mean.

We often prefer a pseudo-random number generator because our simulations results can be exactly reproduced by controlling the seed.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

We often prefer a pseudo-random number generator because
our simulations results can be exactly reproduced by
controlling the seed.
**True.** This is an essential aspect of reproducible data analyses
and simulation studies.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Suppose we have a Pandas Series called **thePop** which contains a census of **25000 subjects**. We also have a simple random sample of **400 individuals** saved in the Series **theSample**. We are interested in studying the behavior of the bootstrap procedure on the simple random sample. Fill in the blanks in the code below to construct **10000 bootstrapped estimates** for the **median**.

```
boot_stats = [
    _____
    .sample(n = ___, replace = __)
    ._____()
    for j in range(_____)
]
```

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Suppose we have a Pandas Series called **thePop** which contains a census of **25000 subjects**. We also have a simple random sample of **400 individuals** saved in the Series **theSample**. We are interested in studying the behavior of the bootstrap procedure on the simple random sample. Fill in the blanks in the code below to construct **10000 bootstrapped estimates** for the **median**.

```
boot_stats = [
      theSample
      .sample(n = 400, replace = True)
      .median()
      for j in range(10000)
      ]
```

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

Consider the following layout of the files A and B onto a distributed file-system of 6 machines.



Assume that all blocks have the same file size and computation takes the same amount of time.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data



| File A | | File B | |
|---|---|---|---|
| A.1 | A.2 | B.1 | B.2 |
| A.3 | A.4 | B.3 | B.4 |

Machine 1

| A.1 | B.2 |
|---|---|
| B.1 | A.4 |

Machine 2

| A.3 | B.1 |
|---|---|
| A.4 | A.2 |

Machine 3

| A.1 | B.4 |
|---|---|
| B.3 | A.3 |

Machine 4

| A.1 | B.2 |
|---|---|
| A.3 | B.3 |

Machine 5

| B.1 | B.3 |
|---|---|
| A.2 | B.4 |

Machine 6

| A.2 | A.4 |
|---|---|
| B.2 | B.4 |

If we wanted to load file A in parallel which of the following
sets of machines would give the best load performance:

1. $\{M1, M2\}$
2. $\{M1, M2, M3\}$
3. $\{M2, M4, M5, M6\}$

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data



| File A | |
|---|---|
| A.1 | A.2 |
| A.3 | A.4 |

| File B | |
|---|---|
| B.1 | B.2 |
| B.3 | B.4 |

Machine 1

| A.1 | B.2 |
| B.1 | A.4 |

Machine 2

| A.3 | B.1 |
| A.4 | A.2 |

Machine 3

| A.1 | B.4 |
| B.3 | A.3 |

Machine 4

| A.1 | B.2 |
| A.3 | B.3 |

Machine 5

| B.1 | B.3 |
| A.2 | B.4 |

Machine 6

| A.2 | A.4 |
| B.2 | B.4 |

If we wanted to load file A in parallel which of the following
sets of machines would give the best load performance:

3  $\{M2, M4, M5, M6\}$

While all choices would be able to load the file, only
$\{M2, M4, M5, M6\}$ could load the file in parallel.
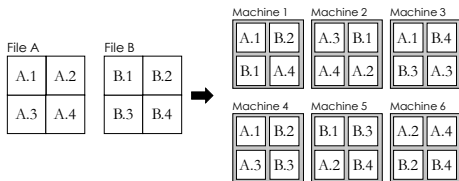
Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

| File A | | File B | |
|---|---|---|---|
| A.1 | A.2 | B.1 | B.2 |
| A.3 | A.4 | B.3 | B.4 |

Machine 1
| A.1 | B.2 |
|---|---|
| B.1 | A.4 |

Machine 2
| A.3 | B.1 |
|---|---|
| A.4 | A.2 |

Machine 3
| A.1 | B.4 |
|---|---|
| B.3 | A.3 |

Machine 4
| A.1 | B.2 |
|---|---|
| A.3 | B.3 |

Machine 5
| B.1 | B.3 |
|---|---|
| A.2 | B.4 |

Machine 6
| A.2 | A.4 |
|---|---|
| B.2 | B.4 |

If we were to lose machines $M1$, $M2$, and $M3$ which of the
following file or files would we lose (select all that apply).

1. File A
2. File B
3. We would still be able to load both files.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data



File A

| A.1 | A.2 |
| A.3 | A.4 |

File B

| B.1 | B.2 |
| B.3 | B.4 |

Machine 1

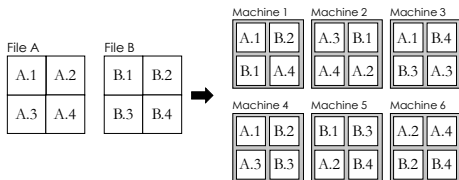| A.1 | B.2 |
| B.1 | A.4 |

Machine 2

| A.3 | B.1 |
| A.4 | A.2 |

Machine 3

| A.1 | B.4 |
| B.3 | A.3 |

Machine 4

| A.1 | B.2 |
| A.3 | B.3 |

Machine 5

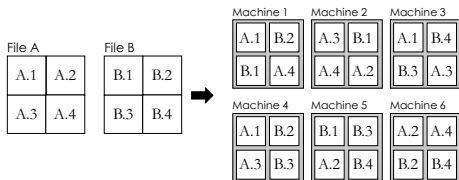| B.1 | B.3 |
| A.2 | B.4 |

Machine 6

| A.2 | A.4 |
| B.2 | B.4 |

If we were to lose machines $M1$, $M2$, and $M3$ which of the
following file or files would we lose (select all that apply).

❸ We would still be able to load both files.

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

| File A | |
|---|---|
| A.1 | A.2 |
| A.3 | A.4 |

| File B | |
|---|---|
| B.1 | B.2 |
| B.3 | B.4 |

Machine 1
| A.1 | B.2 |
|---|---|
| B.1 | A.4 |

Machine 2
| A.3 | B.1 |
|---|---|
| A.4 | A.2 |

Machine 3
| A.1 | B.4 |
|---|---|
| B.3 | A.3 |

Machine 4
| A.1 | B.2 |
|---|---|
| A.3 | B.3 |

Machine 5
| B.1 | B.3 |
|---|---|
| A.2 | B.4 |

Machine 6
| A.2 | A.4 |
|---|---|
| B.2 | B.4 |

If each of the six machines fail with probability $p$, what is the
probability that we will lose block $B$.1 of file B.?

1 $3p$

2 $p^3$

3 $(1 - p)^3$

4 $1 - p^3$

Discussion 12

Exam Review

Probability &
Sampling

EDA &
Visualization

Prediction

Optimization

Inference

Big Data

| File A | |
|---|---|
| A.1 | A.2 |
| A.3 | A.4 |

| File B | |
|---|---|
| B.1 | B.2 |
| B.3 | B.4 |

Machine 1
| A.1 | B.2 |
| B.1 | A.4 |

Machine 2
| A.3 | B.1 |
| A.4 | A.2 |

Machine 3
| A.1 | B.4 |
| B.3 | A.3 |

Machine 4
| A.1 | B.2 |
| A.3 | B.3 |

Machine 5
| B.1 | B.3 |
| A.2 | B.4 |

Machine 6
| A.2 | A.4 |
| B.2 | B.4 |

If each of the six machines fail with probability *p*, what is the probability that we will lose block *B*.1 of file B.?

②  $p^3$