

# CA – Exam 3



## Section Review

**1 Probability Models**

8/15  53%  1:55:34  33%

**2 Statistics**

2/8  25%  52:28  18%

**3 Extended Linear Models**

13/22  59%  1:15:56  49%



1

/1



72%



1.1



2.3



3:11



2:42

**1**

You are given the following information:

- The amount of time one spends in an IRS office is exponentially distributed with a mean of 30 minutes.
- $P_1$  is the probability that John will spend more than an hour of total time in the IRS office, given that he has already been in the office for 20 minutes.
- $P_2$  is the probability that Lucy will spend more than an hour of total time in the IRS office, given that she has just arrived.

**2****3****4****5****6****7****8****9****10****11****12****13****14**

Calculate the difference ( $P_1 - P_2$ ).

8%  A Less than 0.00

16%  B At least 0.00, but less than 0.05

2%  C At least 0.05, but less than 0.10

72%  D At least 0.10, but less than 0.15

2%  E At least 0.15



1/1



73%



1.1



2.5



2:38



3:37

1

You are given the following information:

2

- At the time of installation a given machine has the following survival function:

3

$$S(x) = \left(1 - \frac{x}{100}\right)^{1/2}, \quad \text{for } 0 \leq x \leq 100, \quad x \text{ in months}$$

4

- After 16 months, that machine is still functioning.

5

Calculate the probability that the machine will stop working between months 36 and 51.

6

5%  A Less than 0.095

7

15%  B At least 0.095, but less than 0.105

8

73%  C At least 0.105, but less than 0.115

9

2%  D At least 0.115, but less than 0.125

10

5%  E At least 0.125

11

12

13

14



1

You are given the probability density function for a random variable  $X$ :

2

$$f_X(x) = \left(\frac{\lambda}{\alpha}\right) \left(\frac{x}{\alpha}\right)^{\lambda-1} \exp\left[-\left(\frac{x}{\alpha}\right)^\lambda\right], \quad x \geq 0, \quad \lambda > 0, \quad \alpha > 0$$

3

4

Determine which of the following statements is/are true.

5

6

7

8

9

10

11

12

13

14

I only

B II only

C III only

D I, II, and III

E The correct answer is not given by (A), (B), (C), or (D).



58%



1.3



6.6



13:16



5:41

1

A random variable,  $X$ , is uniformly distributed on the interval  $(m, n)$ .

2

Calculate the  $\text{CTE}_q$  of  $X$ .

3

\* Incorrect Answer

4

A  $\frac{q(n+m)}{2}$

5

$\frac{q(n-m)}{2}$

6

C  $\frac{(1-q)(n+m)}{2}$

7

$\frac{n(1+q) + m(1-q)}{2}$

8

E  $\frac{n(1-q) + m(1+q)}{2}$

9

10

11

12

13

14



1/1

58%

1.4

5.8

1:34

6:22

1

You are given:

2

- Company CA has one representative who answers forum questions.
- Forum questions are posted at a Poisson rate of 10 per hour.
- The time taken for the representative to answer each question is exponentially distributed with a mean of 10 minutes.
- The times questions are posted and the time taken to answer them are independent.
- There are at least two unanswered questions in the system.

3

4

5

Calculate the probability of the representative answering two questions before four new questions are posted.

6

A 11% Less than 0.60

7

58% At least 0.60, but less than 0.62

8

C 14% At least 0.62, but less than 0.64

9

D 6% At least 0.64, but less than 0.66

10

E 11% At least 0.66

11

12

13

14



0/1

71%

1.5

2.9

9:29

4:45

1

You are given the following information about a system of four components:

- The minimal cut sets are  $C_1 = \{1, 4\}$  and  $C_2 = \{2, 3\}$ .
- All components in the system are independent.
- Components 1 and 3 have a reliability of 0.90.
- Components 2 and 4 have a reliability of 0.70.

4

What is the reliability of the system?

5

Incorrect Answer

6

17% A Less than 0.93

7

4% At least 0.93, but less than 0.94

8

71% At least 0.94, but less than 0.95

9

2% D At least 0.95, but less than 0.96

10

6% E At least 0.96

11

12

13

14

 1/1

65%

1.5

4.4

7:24

4:54

1

You are given the following information:

2

- Consider a system of three independent components, each of which functions for an amount of time (in months) uniformly distributed over  $(0, 1)$ .
- Under current design, the system will fail if any of the components fail.
- A new design was made such that the system will fail if 2 or more components fail.

3

Calculate the increase in expected system life in months under the new system design.

4

5

A 8% Less than 0.195

6

B 5% At least 0.195, but less than 0.215

7

C 10% At least 0.215, but less than 0.235

8

D 65% At least 0.235, but less than 0.255

9

E 11% At least 0.255

10

11

12

13

14

 0/1

72%

1.5

4.7

7:19

4:11

1

You are given the following information about a system with three components:

- Two out of three components are required for the system to function.
- Each component has a lifetime that is exponential with mean 1.

Determine which of the following is an expression for the expected lifetime of the system.

**✗ Incorrect Answer**

3% A  $\int_0^{\infty} e^{-3t} dt$

2%  $\int_0^{\infty} te^{-3t} dt$

11% C  $\int_0^{\infty} (e^{-2t} - e^{-3t}) dt$

72% D  $\int_0^{\infty} (3e^{-2t} - 2e^{-3t}) dt$

12% E  $\int_0^{\infty} (3te^{-2t} - 2te^{-3t}) dt$

2

3

4

5

6

7

8

9

10

11

12

13

14



0/1

68%

1.6

4.1

4:30

2:48

1

You are given the following Markov chain transition matrix (states 1, 2, 3, 4, 5):

2

$$\mathbf{P} = \begin{bmatrix} 0.1 & 0.9 & 0.0 & 0.0 & 0.0 \\ 0.5 & 0.5 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.8 & 0.2 \\ 0.0 & 0.0 & 0.6 & 0.4 & 0.0 \\ 0.3 & 0.0 & 0.6 & 0.0 & 0.1 \end{bmatrix}$$

3

4

5

6

7

8

9

10

11

12

13

14

You are also given the following statements:

1. There are two classes in this Markov chain.
2. State 1 is recurrent.
3. State 5 is transient.

Determine which of the following is/are correct.

Incorrect Answer

5%

A

I only

8%

B

II only

8%

C

III only

68%

I, II, and III

12%

The correct answer is not given by (A), (B), (C), or (D).

 1/1

69%



1.6



3.3



8:28



6:33

1

You are given the following information about a homogeneous Markov chain:

- There are three daily wildfire risk states: Green (state 0), Yellow (state 1) and Red (state 2).
- Transition between states occurs at the end of each day.
- The daily transition matrix,  $\mathbf{P} = \begin{bmatrix} 0.82 & m & n \\ 0.61 & 0.28 & 0.11 \\ 0.40 & 0.31 & 0.29 \end{bmatrix}$ .
- The wildfire risk was Yellow on Wednesday.
- Today is Thursday and the wildfire risk is Green.
- The probability that the wildfire risk will be Red on Saturday is 0.07.

2

3

4

5

6

7

8

9

10

11

12

13

14

Calculate the absolute difference between  $m$  and  $n$ .

13%

 A Less than 0.065

6%

 B At least 0.065, but less than 0.075

69%

 C At least 0.075, but less than 0.085

6%

 D At least 0.085, but less than 0.095

6%

 E At least 0.095

 0 / 1

62%



1.6



5.6



12:12



5:27

1

You are given:

2

- Mary plays a game repeatedly.
- Each game ends with her either winning or losing.
- Mary's chances of winning her next game depends on the outcome of the prior game.
  - $P[\text{Winning after a win}] = \min(80\%, P[\text{Winning prior game}] + 10\%)$
  - $P[\text{Winning after a loss}] = 40\%$
- Mary just played her 10<sup>th</sup> game and lost.

3

4

5

6

7

8

9

10

11

12

13

14

Calculate the probability that Mary will lose her 13<sup>th</sup> game.

Incorrect Answer

8%

7%

12%

62%

11%

A

B

C

D

E

Less than 40%

At least 40%, but less than 45%

At least 45%, but less than 50%

At least 50%, but less than 55%

At least 55%





0 / 1



55%



1.7



5.9



12:41



7:17

1

You are given the following information:

- A life insurance company issues a special 3-year discrete insurance to a life  $(x)$ .
- If the policyholder dies, at the end of the year of death, there is a random drawing. With probability 0.2, the death benefit is 50,000. With probability 0.8, the death benefit is 0.
- At the beginning of each year the policy is in effect which  $(x)$  is alive, there is a random drawing. With probability 0.8, the premium  $\pi$  is paid. With probability 0.2, no premium is paid.
- The random drawings are independent.
- The probability of surviving an additional  $k$  years is  ${}_k p_x = 0.9^k$ , for  $k = 0, 1, 2, \dots$
- $i = 0.06$

2

3

4

5

6

7

8

9

10

11

12

13

14

Calculate  $\pi$  using the equivalence principle.

Incorrect Answer

10%

A

Less than 1,100

55%

At least 1,100, but less than 1,200

10%

At least 1,200, but less than 1,300

7%

D

At least 1,300, but less than 1,400

18%

E

At least 1,400



1/1



64%



1.7



4.6



4:47



4:48

1

You are given the following information:

- An annuity-due is issued to a 30-year-old that will pay 1 each year until either she dies or reaches age 50, whichever comes first.
- Mortality follows the Illustrative Life Table.
- Annual interest rate  $i = 6\%$ .

4

Calculate the actuarial present value of this annuity.

5

6%  A Less than 10.0

6

7%  B At least 10.0, but less than 10.5

7

12%  C At least 10.5, but less than 11.0

8

10%  D At least 11.0, but less than 11.5

9

64%  E At least 11.5

10

11

12

13

14



0 / 1



60%



1.8



5.4



5:57



4:20

1

You are given the following simulation process to generate random variable  $X$  using the rejection method:

2

- $X$  has density function:  $f(x) = \frac{2}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right)$ , for  $x > 0$ .

3

- The rejection method is based on  $g(x) = \exp(-x)$ , for  $x > 0$ .

4

- The rejection procedure is as follows:

5

- Step 1: Generate independent random numbers  $Y$  and  $U$ , where  $Y$  follows density function  $g$  and  $U$  is uniform on  $(0,1)$ .

6

- Step 2: If  $U \leq \frac{f(Y)}{cg(Y)}$  stop and set  $X = Y$ . Otherwise return to Step 1.

7

Calculate the minimum possible value for  $c$ .

8

Incorrect Answer

9

5% A Less than 0.50

10

6% At least 0.50, but less than 0.75

11

15% C At least 0.75, but less than 1.00

12

14% D At least 1.00, but less than 1.25

13

60% At least 1.25

14

15

You are given:

16

- $X$  is a beta random variable with  $a = 2, b = 3, \theta = 10$ .
- $X$  is simulated using the acceptance-rejection method with the following random variable with density function:

17

$$g(y) = \frac{1}{10}, \text{ where } 0 < y < 10$$

18

- The following values from the random variable above are given: 7.56, 2.36, 4.51, 8.27.

19

The uniform random numbers are simulated using a mixed congruential generator

20

$$x_{n+1} = (ax_n + c) \bmod m$$

21

with the following parameters:

22

- $x_0 = 300$  (the initial seed)  
 $a = 61$  (the multiplicative factor)  
 $c = 593$  (the additive term)  
 $m = 1,024$  (the modulus)

23

24

Calculate the average of the simulated values of  $X$ .

25

26

6%  A Less than 3

27

61%  B At least 3, but less than 4

28

18%  C At least 4, but less than 5

29

11%  D At least 5, but less than 6

30

5%  E At least 6



0 / 1

66%

2.1

3.6

3:09

5:30

15

16

17

18

19

20

21

22

23

24

25

26

27

28

You have a random sample of five observations:

2.5 7.5 12.5 16.0 17.5

- The probability density function below is fit to the random sample.

$$f(x) = \frac{2}{\beta^2} x_i e^{-\left(\frac{x_i}{\beta}\right)^2}$$

Calculate the maximum likelihood estimate of  $\beta$ .

Incorrect Answer

9%

A

Less than 10

4%

B

At least 10, but less than 11

13%

At least 11, but less than 12

66%

At least 12, but less than 13

8%

E

At least 13



0/1

57%

2.1

5.7

2:42

7:42

15

16

17

18

19

20

21

22

23

24

25

26

27

28

An actuary observes the following 20 losses for a select insurance policy:

20, 25, 36, 38, 42, 52, 55, 57, 65, 66, 69, 71, 72, 73, 74, 74, 74, 80, 81, 82

She believes the distribution which fits the data the best is loglogistic with the following probability density:

$$f(x; \gamma, \theta) = \frac{\gamma x^{\gamma-1}}{\theta^\gamma [1 + (\frac{x}{\theta})^\gamma]^2}$$

She uses the 20<sup>th</sup> and 80<sup>th</sup> percentiles to estimate the two parameters of this distribution.

Calculate the estimated value of  $\theta$ .

Incorrect Answer

3%

A

Less than 20

12%

At least 20, but less than 30

20%

C

At least 30, but less than 40

8%

D

At least 40, but less than 50

57%

At least 50

 0 / 1

72%

2.2

4.0

18:19

5:13

15

16

17

18

19

20

21

22

23

24

25

26

27

28

Determine the Fisher Information of  $n$  independent samples from a geometric distribution with mean  $\beta$ .

**✗ Incorrect Answer**

1% A  $\frac{1}{\beta}$

5%  $\frac{n}{\beta}$

72%  $\frac{n}{\beta} - \frac{n}{\beta + 1}$

9% D  $\frac{n\beta}{\beta + 1}$

13% E  $\frac{\beta(\beta + 1)}{n}$

 0 / 1

64%

2.2

5.0

1:18

6:36

15

16

17

18

19

20

21

22

23

24

25

26

27

28

Let  $X_1, X_2, \dots, X_n$  be a random sample from a population  $X$  with probability mass function:

$$p(x) = \theta(1 - \theta)^x \quad \text{for } x = 0, 1, 2, \dots$$

$\theta$  is an unknown parameter between 0 and 1 and the expected value of  $X$  is:

$$\mathbb{E}(X) = \frac{1 - \theta}{\theta}$$

Determine the Cramer-Rao lower bound for the variance of all unbiased estimators of  $\theta$ .

Incorrect Answer

14% A  $\frac{1}{n}\theta(1 - \theta)$

6%  $n\theta(1 - \theta)$

9% C  $\frac{1}{n^2}\theta(1 - \theta)$

64%  $\frac{1}{n}\theta^2(1 - \theta)$

6% E  $\frac{1}{n^2}\theta^2(1 - \theta)$



0/1



66%



2.3



5.1



4:21



4:18

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

A random sample of 10 screw lengths is taken. The sample mean is 2.5 and the unbiased sample variance is 3.0. The underlying distribution is assumed to be normal.

You want to perform a test of the variance of screw lengths,  $\sigma^2$ .

- Null hypothesis is  $H_0 : \sigma^2 = 6$
- Alternative hypothesis is  $H_1 : \sigma^2 < 6$

Determine the result of this hypothesis test.

Incorrect Answer

11% A Reject  $H_0$  at the 0.005 level

6% Reject  $H_0$  at the 0.010 level, but not at the 0.005 level

9% C Reject  $H_0$  at the 0.025 level, but not at the 0.010 level

8% D Reject  $H_0$  at the 0.050 level, but not at the 0.025 level

66% Do not reject  $H_0$  at the 0.050 level



1/1

61%

2.3

5.1

6:10

6:40

15

16

17

18

19

20

21

22

23

24

25

26

27

28

You are given an independent random sample from a normal distribution  $X$  with unknown mean,  $\mu$ . You are given the following information:

- $\sigma^2(X) = 400$
- $H_0 : \mu = 0$
- $H_1 : \mu = 10$
- You will reject the null hypothesis if  $\bar{X} > \alpha$  for some value of  $\alpha$ .

Calculate the minimum sample size needed so that the probability of a Type I error and the probability of a Type II error are each no more than 10%.

9%

A

Less than 10

10%

B

At least 10, but less than 15

14%

C

At least 15, but less than 20

6%

D

At least 20, but less than 25

61%

At least 25



1/1

78%

2.5

2.9

1:55

4:51

15

16

17

18

19

20

21

22

23

24

25

26

27

28

You are given the following random observations:

0.1 0.2 0.5 1.0 1.3

You test whether the sample comes from a distribution with probability density function:

$$f(x) = \frac{2}{(1+x)^3}, \quad x > 0$$

Calculate the Kolmogorov-Smirnov statistic.

A

Less than 0.18

B

At least 0.18, but less than 0.20

C

At least 0.20, but less than 0.22

D

At least 0.22, but less than 0.24

E

At least 0.24



0 / 1

61%

2.5

5.7

14:34

4:51

15

16

17

18

19

20

21

22

23

24

25

26

27

28

A particular model of car comes in three colors, Blue, Silver, and White. You want to test the null hypothesis that Blue cars are twice as popular as Silver cars, which are twice as popular as White cars.

You take a random sample of 100 cars and collect the following data:

Color	Count
Blue	46
Silver	32
White	22

Calculate the  $p$ -value for this hypothesis test.

Incorrect Answer

6%



Less than 0.005

11%



At least 0.005, but less than 0.010

13%



At least 0.010, but less than 0.025

61%



At least 0.025, but less than 0.050

9%



At least 0.050



0 / 1



59%



3.1



5.5



:59



1:16

15

16

17

18

19

20

21

22

23

24

25

26

27

28

For any statistical learning method, which of the following increases monotonically as flexibility increases?

- I. Training MSE
- II. Test MSE
- III. Bias squared
- IV. Variance

Incorrect Answer

4%

A

I and III

16%

I and IV

3%

C

II and III

18%

D

II and IV

59%

The correct answer is not given by (A), (B), (C), or (D).



1

/1



60%



3.2



6.1



2:33



2:49

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

You are given a random sample of four observations from a population:

$$\{1, 3, 10, 40\}$$

and would like to calculate the standard error of an estimate of the mean of the population using a bootstrap procedure. Your calculation is based on only two bootstrapped data sets.

Calculate the maximum possible value of your standard error estimate.

3% A Less than 10

22% B At least 10, but less than 20

60% C At least 20, but less than 30

10% D At least 30, but less than 40

4% E At least 40



0 / 1



64%



3.2



6.2



2:54



2:50

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

The model  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$  was fit using 6 observations. The estimated parameters are as follows:

- $\hat{\beta}_0 = 2.31$
- $\hat{\beta}_1 = 1.15$
- $s_{\hat{\beta}_0} = 0.057$
- $s_{\hat{\beta}_1} = 0.043$

The following hypothesis test is performed:

- $H_0 : \beta_1 = 1$
- $H_1 : \beta_1 \neq 1$

Calculate the minimum significance level at which the null hypothesis would be rejected.

Incorrect Answer

12% A Less than 0.01

13% At least 0.01, but less than 0.02

64% At least 0.02, but less than 0.05

7% D At least 0.05, but less than 0.10

3% E At least 0.10

15

You are given:

- An actuary wants to study the effect of birth season on the height of individuals.
- Her data set consists of the following 12 individuals:

Individual	Season of Birth	Height	Individual	Season of Birth	Height
1	Spring	155.4	7	Fall	172.7
2	Spring	142.6	8	Fall	165.2
3	Spring	150.6	9	Fall	180.3
4	Summer	155.3	10	Winter	175.6
5	Summer	161.6	11	Winter	172.5
6	Summer	157.0	12	Winter	172.0

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

- She fits an ordinary least squares regression model using Summer as the reference category.
- Below is the partial R output for her regression model:

Coefficients	Estimate
(Intercept)	$b_0$
SeasonSpring	-8.433
SeasonFall	$b_2$
SeasonWinter	$b_3$

Determine which of the following statements about the coefficient estimates can be concluded correctly from the given information.

6%  A  $b_0 = 163.4$

5%  B  $b_2 = 9.333$

63%  C  $b_3 = 15.4$

12%  D All (A), (B), and (C) can be concluded correctly from the given information.

14%  E There is not enough information to infer (A), (B), or (C).



0/1

61%

3.4

6.2

4:16

4:19

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

A botanist wants to investigate the effects of pesticides on the number of strawberries grown. He designed an experiment with three strawberry patches (each with ten plants) where one patch serves as a control, and the other two treated with pesticides A and B, respectively.

After six months, the number of strawberries are recorded. Below are the observations:

$$\sum_{i=1}^3 \sum_{j=1}^{10} (y_{ij} - \bar{y}_i)^2 = 534 \quad \sum_{i=1}^3 \sum_{j=1}^{10} (y_{ij} - \bar{y})^2 = 742$$

where:

- $y_{ij}$  is the number of strawberries grown by plant  $j$  in patch  $i$ .
- $\bar{y}_i$  is the mean number of strawberries grown in patch  $i$ .
- $\bar{y}$  is the overall mean number of strawberries.

The botanist wants to test if pesticides have any impact on the growth of strawberries. He formulates the following hypothesis:

- $H_0$  : The pesticides do not have a significant effect on the growth of strawberries.
- $H_1$  : The pesticides have a significant effect on the growth of strawberries.

Calculate the  $F$  statistic to test the null hypothesis.

✖ Incorrect Answer

17% Less than 4.0

5% At least 4.0, but less than 4.5

8% At least 4.5, but less than 5.0

61% At least 5.0, but less than 5.5

9% At least 5.5



1/1



70%



3.4



4.2



6:20



3:49

15

In an experiment with three blocks, three treatments, and a total of nine observations you are given the following summary of results:

16

	Sum of Squares	Mean Square
Between Treatments	$\frac{146}{9}$	$\frac{73}{9}$
Between Blocks	$\frac{50}{9}$	$\frac{25}{9}$
Residual	$\frac{10}{9}$	$\frac{5}{18}$
Total	$\frac{206}{9}$	

17

18

19

20

21

22

23

24

25

26

27

28

29

30

- $H_0$  : The treatment means are all equal.
- $H_1$  : The treatment means are not all equal.
- A block is defined to be a group of three observations.
- The results are Normally distributed.

Calculate the smallest  $p$ -value for which  $H_0$  can be rejected.

70%



Less than 1%

7%



B At least 1%, but less than 2.5%

13%



C At least 2.5%, but less than 5%

4%



D At least 5%, but less than 10%

6%



E At least 10%



0 / 1

69%

3.5

4.2

3:36

2:58

29

30

Two different data sets were used to construct the four regression models below. The following output was produced from the models:

31

Data Set	Model	Dependent variable	Independent variables	Total Sum of Squares	Residual Sum of Squares
A	1	$Y_A$	$X_{A1}, X_{A2}$	35,930	2,823
A	2	$X_{A1}$	$X_{A2}$	92,990	7,070
B	3	$Y_B$	$X_{B1}, X_{B2}$	27,570	13,240
B	4	$X_{B1}$	$X_{B2}$	87,020	34,650

32

33

34

Determine which one of the following statements best describe the data.

35

Incorrect Answer

36

11% A Collinearity is present in both data sets A and B.

37

12% B Collinearity is present in neither data set A nor B.

38

69% C Collinearity is present in data set A only.

39

6% D Collinearity is present in data set B only.

40

3% E The degree of collinearity cannot be determined from the information given.

41

42

43

44

45



0 / 1



73%



3.5



3.1



2:11



4:00

29

30

You are given:

31

- A linear model,  $Y = \beta_0 + \beta_1 x + \varepsilon$ , is fitted to a data set.
- The data set has 250 observations.

32

- $\sum_{i=1}^{250} x_i = 87,391$
- $\sum_{i=1}^{250} x_i^2 = 36,268,085$

33

Calculate the leverage statistic for an observation with  $x = 1,500$ .

34

Incorrect Answer

35

15% A Less than 0.10

36

6% At least 0.10, but less than 0.20

37

73% At least 0.20, but less than 0.30

38

4% D At least 0.30, but less than 0.40

39

2% E At least 0.40

41

42

43

44

45

23

A normal regression with identity link is performed on a set of data to investigate the average credit card debt of a person. The following tables are results of the regression:

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

Response variable	Average Credit Card Debt	
Response distribution	Normal	
Link	Identity	
Residual Std. Error	102.9	

Parameter	$\hat{\beta}$	Standard Error
Intercept	-547.31	21.4604
Income	-7.80	0.2422
Student		
No	0.00	0.0000
Yes	417.51	17.1716
Rating	3.98	0.05458
Age	-0.62	0.30407

You are also given the following information on one of the observations:

Income	Student	Rating	Age	Average Debt	Leverage
106.25	Yes	483	82	903	0.0377

Calculate the Cook's distance for this observation.

Incorrect Answer

65%



Less than 0.0001

5%



At least 0.0001, but less than 0.0003

8%



C At least 0.0003, but less than 0.0005

5%



D At least 0.0005, but less than 0.0007

17%



E At least 0.0007



23

You are considering using k-fold cross-validation (CV) in order to estimate the test error of a regression model, and have two options for choice of k:

24

- 5-fold CV
- Leave-one-out CV (LOOCV)

25

26

27

28

29

30

31

32

33

34

35

36

4%

A

1-fold CV is usually sufficient for estimating the test error in regression problems.

8%

B

LOOCV and 5-fold CV usually produce similar estimates of test error, so the simpler model is preferable.

8%

C

Running each cross-validation model is computationally expensive.

This is an incorrect answer

67%

Models fit on smaller subsets of the training data result in greater overestimates of the test error.

13%

E

Using nearly-identical training data sets results in highly-correlated test error estimates.



23

You are given the following three statements regarding shrinkage methods in linear regression:

24

I. As tuning parameter,  $\lambda$ , increases towards  $\infty$ , the penalty term has no effect and a ridge regression will result in the unconstrained estimates.

25

II. For a given dataset, the number of variables in a lasso regression model will always be greater than or equal to the number of variables in a ridge regression model.

26

III. The issue of selecting a tuning parameter for a ridge regression can be addressed with cross-validation.

27

6% A I only

28

8% B II only

29

72% C III only

30

8% D I, II and III

31

6% E The answer is not given by (A), (B), (C) or (D).

33

34

35

36



1 / 1

76%

3.8

3.2

1:01

1:10

23

24

25

26

27

28

29

30

31

32

33

34

35

36

You want to investigate whether a particular day of the week affects the number of hourly users of a bike sharing program in your state.

Determine which of the following distribution and link function is most appropriate for this model.

6%



A Normal distribution, identity link function

2%



B Normal distribution, log link function

9%



C Binomial distribution, logit link function

76%



D Poisson distribution, log link function

7%



E Poisson distribution, identity link function



1/1

72%

3.8

3.1

3:43

1:17

32

Andy wants to study the impact of the number of miles driven on the frequency of personal auto insurance claims.

33

Determine which of the following distribution and link function is most appropriate for his model.

34

10% A Normal distribution, identity link function

35

5% B Binomial distribution, logit link function

36

72% C Poisson distribution, log link function

37

8% D Poisson distribution, identity link function

38

5% E Gamma distribution, negative inverse link function

39

40

41

42

43

44

45



1 / 1

64%

3.8

5.0

1:05

1:23

32

Determine which of the following statements are true.

- I. The deviance is useful for testing the significance of explanatory variables in nested models.
- II. The deviance for normal distributions is proportional to the residual sum of squares.
- III. The deviance is defined as a measure of distance between saturated and fitted model.

6% A I only

3% B II only

13% C III only

14% D All but III

64% All

33

34

35

36

37

38

39

40

41

42

43

44

45



1/1



66%



3.8



5.2



5:32



4:22

32

A GLM was used to estimate the expected losses per customer across gender and territory. The following information is provided:

33

- The link function selected is log.
- Q is the base level for Territory.
- Male is the base level for Gender.
- Interaction terms are included in the model.

34

35

36

37

38

39

40

41

42

43

44

45

36

The GLM produced the following predicted values for expected loss per customer:

	Q	R
Male	148	545
Female	446	4,024

Calculate the estimated beta for the interaction of Territory R and Female.

5%

A

Less than 0.85

66%

At least 0.85, but less than 0.95

16%

C

At least 0.95, but less than 1.05

7%

D

At least 1.05, but less than 1.15

6%

E

At least 1.15



32

Three generalized linear models are built to predict the sepal length of irises. Details for these three models are provided below:

33

Model I			Model II			Model III			
Distribution	Normal	Link	Identity	Distribution	Inverse Gaussian <th>Link</th> <td>Inverse-square</td> <th>Distribution</th> <td>Gamma</td>	Link	Inverse-square	Distribution	Gamma
Parameter	Estimate	p-value	Parameter	Estimate	p-value	Parameter	Estimate	p-value	
(Intercept)	1.856	9.85E-12	(Intercept)	0.0639	< 2e-16	(Intercept)	0.2795	< 2e-16	
Sepal.Width	0.6508	< 2e-16	Sepal.Width	-0.0047	6.80E-11	Sepal.Width	-0.0161	5.96E-14	
Petal.Length	0.7091	< 2e-16	Petal.Length	-0.006	< 2e-16	Petal.Length	-0.0189	< 2e-16	
Petal.Width	-0.5565	2.41E-05	Petal.Width	0.0033	4.33E-03	Petal.Width	0.0123	4.01E-04	
Model Statistics			Model Statistics			Model Statistics			
Deviance	14.445		Deviance	0.080		Deviance	0.426		
AIC	84.643		AIC	95.898		AIC	82.629		

34

35

36

37

38

39

40

41

42

43

44

45

The Akaike Information Criterion is used to select the best out of these three models.

Calculate the predicted sepal length for an iris with a sepal width of 4.4, a petal length of 6.9, and a petal width of 2.5.

9.17

(round to the nearest 0.01)

Correct Answer: 9.17



1

/1



54%



3.9



6.0



:39



1:13

32

Determine which of the following link functions are suitable for a generalized linear model used to model the likelihood of a fraudulent claim. (select all that apply)

33

62% Complementary log-log link function

34

14% B Identity link function

35

13% C Inverse link function

36

82% D Logit link function

38

6% E Power link function

39

40

41

42

43

44

45



32

33

34

35

36

37

38

39

40

41

42

43

44

45

A Poisson regression with identity link is performed to predict the annual number of claims. The table below shows the estimated parameters:

<b>Response Variable</b>	Annual Number of Claims
<b>Response Distribution</b>	Poisson
<b>Link</b>	Identity

Parameter	df	$\hat{\beta}$	p-value
Intercept	1	0.712	<0.0001
Age	1	0.003	<0.0001
Gender: Male	1	-1.122	<0.0001
Gender: Female	0	0.000	
Number of Legal Drivers in Household	1	0.465	<0.0001
Zone: A	0	0.000	
Zone: B	1	0.468	<0.0001
Zone: C	1	1.182	<0.0001

Calculate the variance of the annual number of claims for a male age 35 from zone A with 3 legal drivers in his household.

5%



Less than 1.0

75%



At least 1.0, but less than 2.0

17%



At least 2.0, but less than 3.0

2%



At least 3.0, but less than 4.0

2%



At least 4.0



32

You have fit a generalized linear model using a Poisson distribution. One observation of the response variable is 59. The corresponding GLM fitted value is 71. The corresponding leverage of this point is 22%.

33

Calculate the standardized deviance residual corresponding to this observation.

34

Incorrect Answer

35

13% A Less than -1.8

36

7% At least -1.8, but less than -1.7

37

58% At least -1.7, but less than -1.6

38

6% D At least -1.6, but less than -1.5

39

16% E At least -1.5

40

41

42

43

44

45

**32**

Suppose we have a linear regression model

**33**

$$Y = \beta_0 + \beta_1 \cdot b_1(X) + \beta_2 \cdot b_2(X) + \varepsilon$$

**34**

with basis functions  $b_1(X) = X$  and  $b_2(X) = X^2$ . This model can be rewritten as

**35**

$$Y = \alpha_0 + \alpha_1 \cdot b_1^*(X) + \alpha_2 \cdot b_2^*(X) + \varepsilon$$

**36**

where  $b_1^*(X) = X - 1$  and  $b_2^*(X) = (X - 1)^2$ .

**37**

Which of the following is the correct expression for  $\alpha_0$  in terms of the  $\beta$  parameters?

**38**

Incorrect Answer

**39**

3%    A     $\beta_0 + 2\beta_1$

**40**

67%        $\beta_0 + \beta_1 + \beta_2$

**41**

17%    C     $\beta_0 + 2\beta_1 + \beta_2$

**42**

5%    D     $\beta_0 + 2\beta_1 + 2\beta_2$

**43**

8%    The answer is not given by (A), (B), (C) or (D).

**44****45**



1

/1



67%



3.11



4.8



1:16



1:09

32

A modeler creates a local regression model. After reviewing the results, the fitted line appears too wiggly, over-responding to trends in nearby data points. The modeler would like to adjust the model to produce more intuitive results.

33

Determine which one of the following adjustments the modeler should make.

34

A Add a linear constraint in the regions before and after the first knot

35

B Increase the number of orders in the regression equation

36

C Increase the number of knots in the model

37

D Reduce the number of knots in the model

38

✓ Increase the span,  $s$ , of the model

40

41

42

43

44

45

23

You are fitting a linear local regression model

24

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

25

to the following set of 20 data points using a span of  $s = 0.2$ :

26

27

Obs.	1	2	3	4	5	6	7	8	9	10
$y_i$	2.7	3.7	6.4	11.9	16.4	15.4	15.8	16.1	17.0	19.2
$x_i$	2.2	4.8	6.9	9.2	10.2	10.8	12.0	13.8	13.9	14.1

28

29

30

Obs.	11	12	13	14	15	16	17	18	19	20
$y_i$	20.1	21.3	24.7	21.9	26.4	25.5	28.2	29.1	26.7	27.1
$x_i$	15.2	15.4	15.9	16.4	17.5	18.8	19.0	19.4	22.8	26.5

31

32

33

34

35

36

37

38

39

40

41

42

43

When fitting the linear local regression at  $X = 10$ , you are using the weighting function

$$K_i = \frac{2 - |10 - x_i|}{2}$$

This linear local regression at  $X = 10$  has produced the estimate  $\hat{\beta}_1 = 2.489$ .Calculate the estimated value of  $\hat{\beta}_0$ .

Incorrect Answer

8% A Less than -13

9% At least -13, but less than -11

73% At least -11, but less than -9

6% D At least -9, but less than -7

4% E At least -7