

PART I: Multiple Choice (5 Points Per Question). Choose the **best** answer.

1. Which of the following is a linear model?

- (a) $E(Y|X = x) = (\beta_0 + \beta_1 x)^2$
- (b) $E(Y|X = x) = \beta_0 + \beta_1 x + \beta_2 x^2$
- (c) $E(Y|X = x) = \sqrt{\beta_0 + \beta_1 x + \beta_2 x^2}$
- (d) $E(Y|X = x) = 1/(\beta_0 + \beta_1 x)$

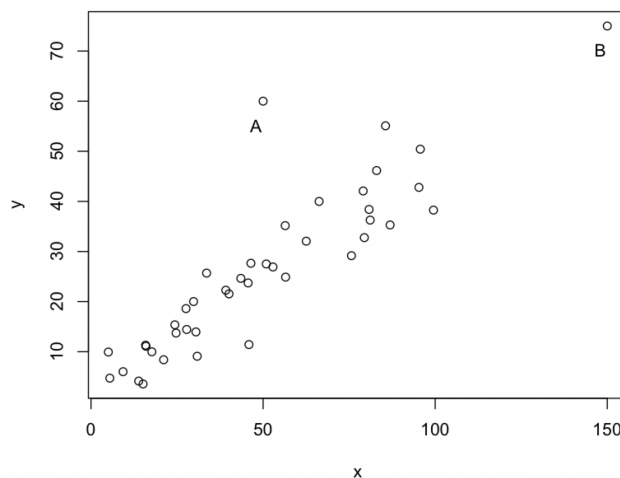
2. Which of the following must be true in order for a simple linear regression model to be valid?

- (a) The data must be collected across time.
- (b) The errors must have constant variance.
- (c) The errors must be normally distributed. (NOT NECESSARILY REQUIRED IF LARGE SAMPLE SIZE)
- (d) The relationship between x and y must be exponential.
- (e) The relationship between x and y must be quadratic.

3. Consider a dataset with response y and predictor variable x . Suppose we fit the model $\sqrt{Y} = \beta_0 + \beta_1 x + e$. A 95% confidence interval for the mean response with $x = 10$ is $(1.25, 4.15)$. The MSE of the model was 0.05. How should the confidence interval be back-transformed?

- (a) $(e^{1.25+0.05/2}, e^{4.15+0.05/2})$
- (b) $(e^{1.25-0.05}, e^{4.15+0.05})$
- (c) $(1.25^2 + 0.05, 4.15^2 + 0.05)$ (SEE SLIDE 80 OF CHAPTER 3 NOTES)
- (d) $(1.25^2 - 0.05, 4.15^2 + 0.05)$
- (e) $(\frac{1}{1.25} (1 + \frac{0.05}{1.25}), \frac{1}{4.15} (1 + \frac{0.05}{4.15}))$

4. The points labeled “A” and “B” in the graph below can be described using which of the following?



- (a) Point A is a bad leverage point.
 - (b) Point A is a good leverage point.
 - (c) Point B is a bad leverage point.
 - (d) **Point B is a good leverage point. (IT IS A LEVERAGE POINT BECAUSE ITS x VALUE IS FAR FROM THE BULK OF THE OTHERS. IT IS A “GOOD” LEVERAGE POINT BECAUSE IT STRENGTHENS THE SLOPE ESTIMATE WE HAD USING THE REST OF THE POINTS.)
5. Which of the following is **not** a reason for transforming predictor and / or response variables?
- (a) To ensure the relationship between the predictor and response is a straight line.
 - (b) **To transform outliers into leverage points.
 - (c) To stabilize the variance of the residuals.
 - (d) To reduce the influence of outliers.
 - (e) To estimate percentage effects (elasticity).
6. What should be done with an outlier?
- (a) Always remove outliers.
 - (b) **Only remove an outlier if there is a plausible way in which it is different from the rest of the observations.
 - (c) Remove an outlier if including it gives unexpected results.
 - (d) Only remove an outlier if it is also a good leverage point.
7. Suppose we have a large sample size and fit a simple linear regression model. Residuals are required to be Normally distributed for which of the following?
- (a) **For prediction intervals for an individual response.
 - (b) For prediction intervals for a mean response.
 - (c) For confidence intervals for an individual response.
 - (d) For confidence intervals for a mean response.

PART II: Short Answer (8 Points Each Part)

8. An advertising executive is interested in the number of clicks on his company’s website within one hour of running television ads. He randomly chose four prime-time one-hour slots and ran either 1, 2, 3, or 4 ads in each hour. There were four responses measured, y_1, y_2, y_3, y_4 , depending on the number of ads run.
- (a) For the model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$, where \mathbf{y} , $\boldsymbol{\beta}$, and \mathbf{e} are vectors and \mathbf{X} is the design matrix, write out \mathbf{X} and $\boldsymbol{\beta}$. Assume there is both an intercept β_0 and a slope β_1 .

$$\mathbf{X} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

- (b) Interpret the model coefficients in the context of the problem (don't just say one is the intercept and the other is the slope).

THE INTERCEPT β_0 IS THE MEAN RESPONSE ASSOCIATED WITH 0 ADS. THIS IS NOT MEANINGFUL FOR US, SINCE WE DO NOT HAVE ANY DATA FOR 0 ADS. IF WE MEAN-CENTERED THE EXPLANATORY VARIABLE, THE INTERCEPT WOULD REPRESENT THE MEAN RESPONSE ASSOCIATED WITH RUNNING 2.5 ADS PER HOUR. THIS AGAIN IS OF LIMITED VALUE, SINCE WE DO NOT HAVE DATA WITH 2.5 ADS.

THE SLOPE β_1 IS THE MEAN CHANGE IN RESPONSE ASSOCIATED WITH INCREASING THE NUMBER OF ADS PER HOUR BY ONE UNIT. FOR EXAMPLE, THE MEAN CHANGE IN NUMBER OF CLICKS COMPARING 2 ADS PER HOUR TO 1 AD PER HOUR IS β_1 . SIMILARLY, THE MEAN CHANGE IN NUMBER OF CLICKS COMPARING 4 ADS PER HOUR TO 1 AD PER HOUR IS $3\beta_1$.

- (c) Use the design matrix above and the general least squares solution for the parameter estimate vector to calculate estimates for the parameters.

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix} = \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}$$

$$(\mathbf{X}'\mathbf{X})^{-1} = \frac{1}{120 - 100} \begin{bmatrix} 30 & -10 \\ -10 & 4 \end{bmatrix} = \frac{1}{20} \begin{bmatrix} 30 & -10 \\ -10 & 4 \end{bmatrix}$$

$$\begin{aligned} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} &= \frac{1}{20} \begin{bmatrix} 30 & -10 \\ -10 & 4 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \\ &= \frac{1}{20} \begin{bmatrix} 20 & 10 & 0 & -10 \\ -6 & -2 & 2 & 6 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \\ &= \frac{1}{20} \begin{bmatrix} 20y_1 + 10y_2 - 10y_4 \\ -6y_1 - 2y_2 + 2y_3 + 6y_4 \end{bmatrix} \end{aligned}$$

9. Recall the 2016 / 2020 presidential election data that we saw in class. The response variable is the percentage of Trump voters in each of the 50 states in 2020, while the explanatory variable is the percentage of Trump voters in each state in 2016. The model summary from R is shown below:

Call:

```
lm(formula = Trump_2020 ~ Trump_2016)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-0.159344	-0.007675	0.008245	0.019979	0.125834

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.05547	0.03605	1.539	0.13
Trump_2016	0.88598	0.07180	12.340	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.05136 on 48 degrees of freedom

Multiple R-squared: 0.7603, Adjusted R-squared: 0.7553

F-statistic: 152.3 on 1 and 48 DF, p-value: < 2.2e-16

For use below, here are a few critical values:

$z_{0.025}$	$t_{0.025,45}$	$t_{0.025,46}$	$t_{0.025,47}$	$t_{0.025,48}$	$t_{0.025,49}$	$t_{0.025,50}$
1.9600	2.0141	2.0129	2.0117	2.0106	2.0096	2.0086

A few more pieces of information that you may need:

- `mean(Trump_2016) = 0.4918`
- `mean(Trump_2020) = 0.4912`
- `sum((Trump_2016 - mean(Trump_2016)) ^ 2) = 0.5117`
- `sum((Trump_2020 - mean(Trump_2020)) ^ 2) = 0.5283`

- (a) Compute a 95% confidence interval for the mean percentage voting Trump in 2020 for states that voted 60% Trump in 2016.

FIRST NOTE THAT THE FITTED VALUE IS

$$\widehat{\text{TRUMP}}_{2020} = 0.0555 + 0.6(0.8860) = 0.5871$$

THE FORMULA FOR A CONFIDENCE INTERVAL ON A MEAN RESPONSE WHEN $x = x^*$ IS

$$\hat{y} \pm t_{\alpha/2, n-2} s \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX}}$$

USING THE PROVIDED INFORMATION, THIS IS

$$0.5871 \pm 2.0106(0.0514) \sqrt{\frac{1}{50} + \frac{(0.6 - 0.4918)^2}{0.5117}} = (0.5657, 0.6085)$$