

FINAL EXAMINATION VERSION A

INSTRUCTIONS: This exam is closed-book, closed-notes. You may use a calculator on this exam, but not a graphing calculator, a calculator with alphabetical keys, nor a mobile phone. Point values for each question are noted in brackets. Tables of the t-distribution, the F-distribution, and the chi-square distribution are attached. Maximum total points are 200.

NOTATION: In this exam, $\hat{\beta}_j$ denotes the least-squares coefficient estimators of the equation $y_i = \beta_1 + \beta_2 x_{i2} + \dots + \beta_K x_{iK} + \varepsilon_i$. The least-squares fitted value is denoted \hat{y}_i . The least-squares residual is denoted $\hat{\varepsilon}_i$. The sample size is denoted n . The true or population value of the variance of the unobserved error term ε_i is denoted σ^2 . The (unbiased) least-squares estimator of σ^2 is denoted $\hat{\sigma}^2$. The sample mean of y is denoted \bar{y} . The natural logarithm is denoted $\ln(\cdot)$.

I. MULTIPLE CHOICE: Circle the one best answer to each question. Use margins for scratch work [2 pts each—44 pts total]

(1) A data set that includes many observations at the same point in time is called

- a. a cross-section.
- b. a time-series.
- c. a pooled data set.
- d. a panel data set.

(2) Suppose we wish to fit the equation $y = \beta_1 + \beta_2 x$ to data by the method of *least absolute deviation*. This method minimizes which function of the data?

- a. $f(\beta_1, \beta_2) = \sum (y_i - \beta_1 - \beta_2 x_i)$.
- b. $f(\beta_1, \beta_2) = \sum (y_i^2 - (\beta_1 + \beta_2 x_i)^2)$.
- c. $f(\beta_1, \beta_2) = \sum (y_i - \beta_1 - \beta_2 x_i)^2$.
- d. $f(\beta_1, \beta_2) = \sum |y_i - \beta_1 - \beta_2 x_i|$.
- e. $f(\beta_1, \beta_2) = \sum (\beta_1 + \beta_2 x_i)^2$.

(3) An estimator $\hat{\theta}$ of an unknown population parameter θ is said to be *asymptotically unbiased* if

- a. $E(\hat{\theta}) = \theta$.
- b. $E(\hat{\theta}) = 0$.
- c. $\lim_{n \rightarrow \infty} E(\hat{\theta}) = 0$.
- d. $\lim_{n \rightarrow \infty} E(\hat{\theta}) = \theta$.
- e. $\lim_{n \rightarrow \infty} \text{Prob}(|\hat{\theta} - \theta| > \delta) = 0$, for all $\delta > 0$.

(4) Suppose the p-value for a test statistic is 0.064. If the size of the test is 5 percent, we

- a. can reject the null hypothesis.
- b. cannot reject the null hypothesis.
- c. cannot compute the test statistic.
- d. answer cannot be determined from the information given.

(5) Suppose the equation $y = \beta_1 + \beta_2 x$ is fitted to n observations on x and y by the method of least squares. Then the equation

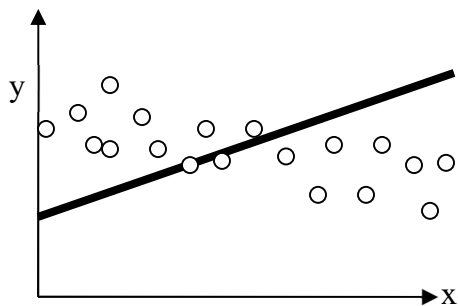
$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n \hat{\varepsilon}_i^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

holds

- if the fit is good.
- if the fit is poor.
- if there are enough observations.
- always.

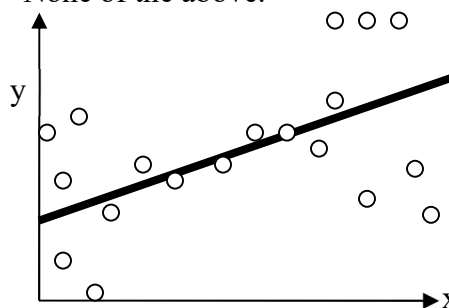
(6) In the graph below, the solid line is the true population regression line and the circles are observations in the sample. Which assumption appears to be violated in this sample?

- $E(\varepsilon_i|x_i) = 0$.
- Homoskedasticity: $\text{Var}(\varepsilon_i) = \sigma^2$, a constant.
- No autocorrelation: $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$ for $i \neq j$.
- All of the above.
- None of the above.



(7) In the graph below, the solid line is the true population regression line and the circles are observations in the sample. Which assumption appears to be violated in this sample?

- $E(\varepsilon_i|x_i) = 0$.
- Homoskedasticity: $\text{Var}(\varepsilon_i) = \sigma^2$, a constant.
- No autocorrelation: $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$ for $i \neq j$.
- All of the above.
- None of the above.



- (8) The LS predictor of \hat{y}_{n+1} given x_{n+1} differs from the actual value y_{n+1} because
- the LS estimates $\hat{\beta}_1$ and $\hat{\beta}_2$ differ from the true parameter values β_1 and β_2 .
 - the actual value y_{n+1} depends in part on a new error term ε_{n+1} .
 - both of the above.
 - They do not differ. The LS predictor equals exactly the actual value y_{n+1} by definition.

(9) Suppose we have the demand function

$$\ln(y) = 15 - 0.03x,$$

where y denotes quantity demanded of gasoline in gallons, and x denotes price per gallon in dollars. Which of the following is true?

- a. If the price increases by 1 percent, then quantity demanded decreases by 3 percent.
- b. If the price increases by 1 dollar, then quantity demanded decreases by 3 percent.
- c. If the price increases by 1 percent, then quantity demanded decreases by 0.03 gallons.
- d. If the price increases by 1 dollar, then quantity demanded decreases by 0.03 gallons.

(10) A *high leverage point* is an observation

- a. whose y -value is quite different from the rest of the sample.
- b. whose x -value is quite different from the rest of the sample.
- c. that contains missing values for x and/or y .
- d. that lies exactly on the true regression line.

(11) *Heteroskedasticity* means that

- a. some x variables are perfectly correlated with each other.
- b. random error terms of different observations are correlated.
- c. random error terms of different observations have different variances.
- d. random error terms are correlated with one or more of the x variables.

(12) Suppose we estimate the equation

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4},$$

and we want to test the null joint hypothesis that $\beta_2 = \beta_3 = \beta_4 = 0$. We should reject the null hypothesis at 5% significance if

- a. *all* the t -statistics for β_2 , β_3 , and β_4 are greater in absolute value than their 5% critical points.
- b. *any* of the t -statistics for β_2 , β_3 , or β_4 are greater in absolute value than their 5% critical points.
- c. THE F statistic is less than its 5% critical point.
- d. THE F statistic is greater than its 5% critical point.
- e. THE F statistic is either less than its lower 2.5 % critical point or greater than its upper 2.5 % critical point.

(13) Suppose we estimate the equation

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3}.$$

If x_{i2} and x_{i3} are *perfectly* correlated, then the least-squares estimators of their coefficients

- a. will have large standard errors.
- b. will be zero.
- c. cannot be computed.
- d. will be biased.
- e. will be inconsistent.

(14) A *dummy variable* means

- a. a regressor that should never been included in the regression.
- b. a regressor that takes only two values, zero and one.
- c. a regressor whose value never changes throughout the sample.
- d. a regressor used in place of another regressor.

(15) Suppose we wish to estimate the relationship between the size of a firm and the number of patents it receives using data on 1000 business firms. Moreover, we want to allow the intercept to be different by industry (manufacturing, hospitality, retail, etc.). If we have five industries in our data, we need

- a. one dummy variable.
- b. two dummy variables.
- c. three dummy variables.
- d. four dummy variables.
- e. five dummy variables.

(16) If the purpose of our regression is *causal inference*, then we should include additional regressors if they

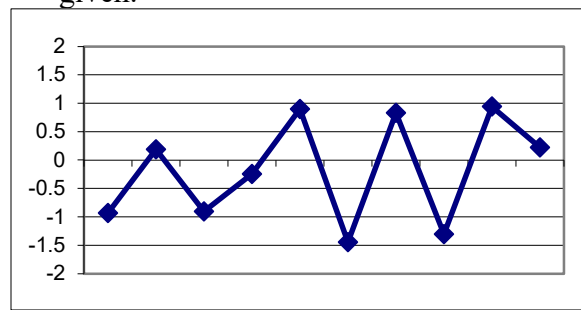
- a. increase the degrees of freedom.
- b. prevent omitted-variable bias.
- c. improve the fit of the equation.
- d. raise the sum of squared residuals.

(17) If the random error term is heteroskedastic then the least squares estimators of the coefficients will

- a. be biased.
- b. be inconsistent.
- c. be impossible to compute.
- d. have incorrect standard errors.

(18) The time series u_t graphed below has mean zero. It appears to be

- a. positively serially-correlated.
- b. negatively serially-correlated.
- c. serially uncorrelated.
- d. Cannot be determined from information given.



(19) If the random process u_t has a unit root, then

- a. $u_t = 1$.
- b. u_t wanders away from its mean.
- c. $E(u_t) = 1$.
- d. $\text{Var}(u_t) = 1$.
- e. the square root of $u_t = 1$.
- f. All of the above.

(20) Suppose u_t is defined by

$$u_t = 3.1 + u_{t-1} + \varepsilon_t,$$

where ε_t is an independent identically-distributed process with mean zero and u_0 is nonrandom. Which is true?

- a. The variance of u_t depends on t .
- b. The mean of u_t depends on t .
- c. The series u_t does not tend to return to a trend line.
- d. All of the above.

(21) To identify and estimate any *integrated time-series process*, we must first compute

- a. logarithms of the data.
- b. first differences of the data.
- c. differences of the data from the sample mean.
- d. cumulative sums of the data.

(22) Given the random walk model with drift

$$y_t = 5 + y_{t-1} + \varepsilon_t,$$

and $y_T = 1.8$, the two-steps-ahead forecast of y_{T+2} is

- a. 0.
- b. 1.8.
- c. 5.0.
- d. 11.8.

II. SHORT ANSWER: Please write your answers in the boxes on this question sheet. Use margins for scratch work.

(1) [Summation operator: 6 pts] Let $\bar{x} = \frac{1}{n} \sum x_i$ and consider the following equations. Write “TRUE” for equations that are necessarily true for any values of x_i and y_i . Write “FALSE” for equations that are not necessarily true. All sums run from $i = 1, \dots, n$.

a. $\sum (x_i - \bar{x}) = 0$

b. $\sum x_i^2 = (\sum x_i)^2$

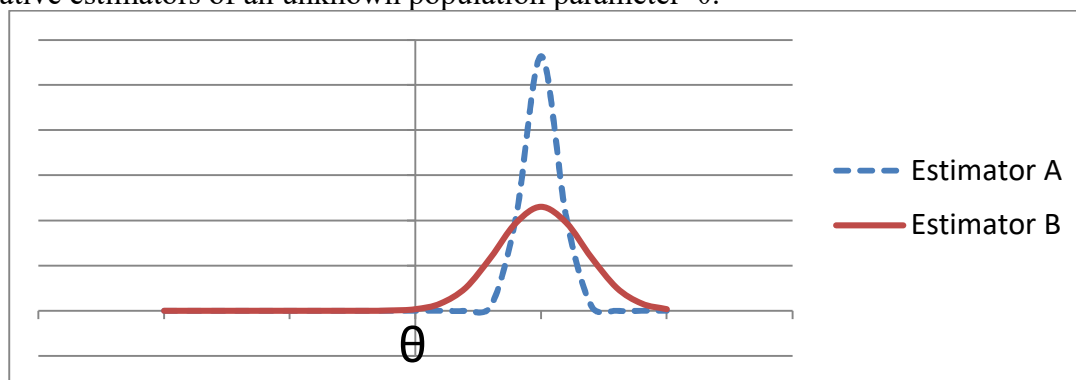
c. $\sum (\bar{x} x_i) = n \bar{x}^2$

(2) [Mean and variance of linear function: 4 pts] Suppose X is a random variable with mean $E(X) = 2$ and variance $\text{Var}(X) = 3$. Suppose $Y = 1 + 2X$.

a. Compute the mean of Y , that is, $E(Y)$.

b. Compute the variance of Y , that is, $\text{Var}(Y)$.

(3) [Properties of estimators: 4 pts] The graph below shows the density functions for two alternative estimators of an unknown population parameter θ .



a. Which estimator has greater **bias**? Answer “A,” “B,” or “EQUAL.”

b. Which estimator has greater **variance**? Answer “A,” “B,” or “EQUAL.”

(4) [Algebraic properties: 8 pts] Suppose the equation $y_i = \beta_1 + \beta_2 x_i + \varepsilon_i$ is fitted by least squares. Which equations hold necessarily, regardless of the data? Write "TRUE" or "FALSE" in the boxes below.

a. $\sum y_i \hat{y}_i = 0$

--

b. $\sum \hat{\varepsilon}_i = 0$

--

c. $\sum x_i \hat{\varepsilon}_i = 0$

--

d. $\sum (x_i - \bar{x}) = 0$

--

(5) [Properties: 10 pts] Which assumptions are required for the least-squares estimators to be *unbiased* estimators? Write "REQUIRED" or "NOT REQUIRED" in the boxes below.

a. Variance of error term is zero: $\text{Var}(\varepsilon_i) = 0$.

--

b. Error term is normally-distributed: $\varepsilon_i \sim N(0, \sigma^2)$

--

c. Homoskedasticity: $\text{Var}(\varepsilon_i) = \sigma^2$, a constant.

--

d. Conditional mean of error term is zero: $E(\varepsilon_i|x_i) = 0$.

--

e. No autocorrelation: $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$ for $i \neq j$.

--

(6) [Properties: 10 pts] Which assumptions are required for the least-squares estimators to be *best linear unbiased estimators* (BLUE)? Write "REQUIRED" or "NOT REQUIRED" in the boxes below.

a. Variance of error term is exactly one: $\text{Var}(\varepsilon_i) = 1$.

--

d. Conditional mean of error term is zero: $E(\varepsilon_i|x_i) = 0$.

--

e. No autocorrelation: $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$ for $i \neq j$.

--

b. Error term is normally-distributed: $\varepsilon_i \sim N(0, \sigma^2)$

--

c. Homoskedasticity: $\text{Var}(\varepsilon_i) = \sigma^2$, a constant.

--

(7) [Variance of LS estimators: 8 pts] Suppose we estimate the equation

$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i$. Assume that the Gauss-Markov assumptions hold. Answer TRUE or FALSE below: The variance of the least-squares slope estimator $\hat{\beta}_2$ is smaller, and thus the true value of β_2 is estimated more precisely,

- the closer the correlation between x_{i2} and x_{i3} .
- greater the variation of x_{i2} around the sample mean \bar{x}_2 .
- the larger the variance of the error term $\text{Var}(\varepsilon_i) = \sigma^2$.
- the larger the sample size.

(8) [Adding regressors: 10 pts] Suppose we first estimate the equation $y_i = \beta_1 + \beta_2 x_{i2}$ by least squares, and then estimate $y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3}$. What are the consequences of adding the regressor x_{i3} ? In each box below, write one of the following:

- “must increase,”
- “must decrease,”
- “can either increase or decrease,”
- “must remain constant.”

- The standard errors of the estimated coefficients...
- The sum of squared residuals...
- The R^2 value...
- Theil’s adjusted R^2 (also called “ \bar{R}^2 ”)
- The values of the estimated coefficients $\hat{\beta}_1$ and $\hat{\beta}_2$...

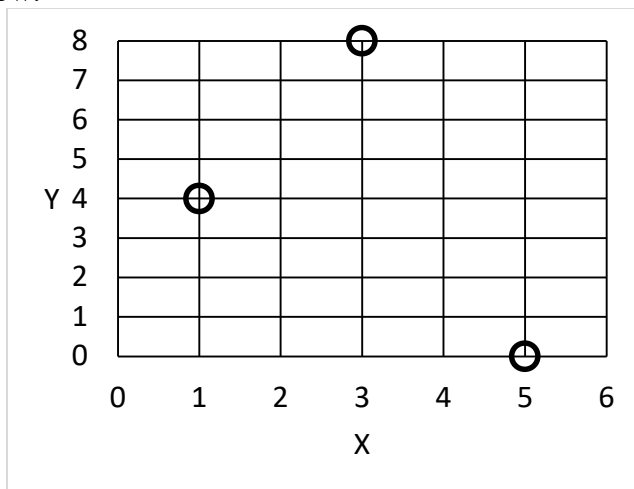
(9) [Stochastic processes: 10 pts] Let ε_t denote an independent identically-distributed (IID) process, and consider the stochastic processes defined by the equations below. Identify each process by filling in three numbers for ARIMA(p, d, q).

- $u_t = 0.4 u_{t-1} + \varepsilon_t$.
- $u_t = u_{t-1} + \varepsilon_t$.
- $u_t = \varepsilon_t + 0.6 \varepsilon_{t-1} - 0.1 \varepsilon_{t-2}$.
- $u_t = 0.6 u_{t-1} + \varepsilon_t + 0.2 \varepsilon_{t-1}$.
- $u_t = 0.5 u_{t-1} + 0.2 u_{t-2} + \varepsilon_t + 0.3 \varepsilon_{t-1}$.

ARIMA(, ,)
ARIMA(, ,)
ARIMA(, ,)
ARIMA(, ,)
ARIMA(, ,)

III. PROBLEMS: Please write your answers in the boxes on this question sheet. Show your work and circle your final answers.

(1) [Definition of least-squares: 15 pts] Suppose we have three observations on X and Y shown in the graph below.



It can be shown that the least-squares estimate for the intercept is $\beta_1 = 7$ and the least squares estimate for the slope is $\beta_2 = -1$.

- a. Compute the three fitted values \hat{y}_i of this least-squares estimated regression line.

- b. Sketch the least-squares fitted line in the graph above.
c. Compute the sample means (\bar{x}, \bar{y}) and verify numerically that the fitted line passes through the sample means.

- d. Compute the three residuals $\hat{\varepsilon}_i$ of this estimated least-squares regression line.

- e. Compute the sum of squared residuals.

(2) [LS confidence intervals, tests, elasticity: 21 pts] Suppose we estimate the effect of temperature on ice cream sales using a sample of $n=300$ days. Let y_i denote daily ice cream sales and x_i denote temperature. The model $y_i = \beta_1 + \beta_2 x_i$ is estimated with the following results. Numbers in parentheses are standard errors.

Ice cream	=	65.0	+	22.0	Temperature
sales		(20.0)		(5.0)	

- a. [3 pts] Suppose the temperature is 70 degrees. According to these results, what are predicted sales?

- b. [3 pts] Suppose the temperature increases by 10 degrees. By how much would ice cream sales increase? That is, what is the predicted change Δy when $\Delta x = 10$?

- c. [6 pts] Compute a **95%** confidence interval for the **intercept, β_1** .

- d. [9 pts] Test the hypothesis that temperature has a **positive** effect on ice cream sales, against the null hypothesis that temperature has no effect (a **one-tailed test**) at **5%** significance. Give the value of the test statistic, the critical point(s) from a table, and your conclusion (whether you can reject null hypothesis).

Value of test statistic = _____. Critical point = _____.

Can you reject null hypothesis? _____.

(3) [Dummy variables and structural change: 20 pts] Suppose we wish to estimate the effect of income on spending for travel, using a sample of **120** households.

spending = spending for travel by household.
 income = income of household.
 retired = 1 if all members of household are retired.
 = 0 if at least one member is still working.

The following four equations were estimated, with the sums of squared residuals (SSR) as shown.

- | | | |
|-----|--|---------|
| [1] | $\text{spending}_i = 56 + 0.2 \text{ income}$ | SSR=372 |
| [2] | $\text{spending} = 46 + 0.3 \text{ income} + 8.0 \text{ retired}$ | SSR=351 |
| [3] | $\text{spending} = 52 + 0.5 \text{ income} - 0.1 (\text{retired} \times \text{income})$ | SSR=350 |
| [4] | $\text{spending} = 33 + 0.4 \text{ income} + 9.0 \text{ retired}$
$\quad - 0.1 (\text{retired} \times \text{income})$ | SSR=232 |

First, consider equation [4].

- According to equation [4], what is the intercept for non-retired households?
- According to equation [4], what is the intercept for retired households?
- According to equation [4], what is the slope for retired households?

Second, test the null hypothesis that all households have the same intercept and slope, against the alternative hypothesis that the intercept for retired households is different from the intercept for nonretired households (but the slope is the same) at 5% significance. Assume the Gauss-Markov assumptions are satisfied and the error term is normally-distributed.

- Which equation, [1], [2], [3], or [4], is the *restricted* equation, representing the null hypothesis?
- Which equation, [1], [2], [3], or [4], is the *unrestricted* equation, representing the alternative hypothesis?
- [10 pts] Give the value of the test statistic, its degrees of freedom, the critical point, and your conclusion (whether you can reject the null hypothesis).

Degrees of freedom in numerator = _____ Degrees of freedom in denominator = _____

Value of F statistic = _____ Critical point = _____

Reject null hypothesis? _____.

(4) [Forecasting, trends and seasonality: 8 pts] Suppose we have estimated the following model for sales, using 80 quarterly observations from the first quarter of 2005 to the fourth quarter of 2024.

$$sales_t = 8.3 + 0.2 \text{ trend} - 0.6 q1_t - 0.2 q2_t - 0.5 q3_t + \varepsilon_t.$$

The regressor “trend” equals 1 in the first quarter of 2005, equals 2 in the second quarter of 2005, and so forth, and equals 80 in the fourth quarter of 2024. The regressors “q1,” “q2,” and “q3,” are quarterly dummy variables for the first, second and third quarters respectively. The error term ε_t is an independent, identically-distributed process with $E(\varepsilon_t) = 0$ and $\text{Var}(\varepsilon_t) = \sigma^2$, constant.

- a. [2 pts] If a dummy variable “q4” for the fourth quarter were also included, then what econometric problem would result?

- b. Compute the forecast of sales in the first quarter of 2025.

- c. Compute the forecast of sales in the second quarter of 2025.

(5) [Forecasting, AR model: 9 pts] Suppose we have estimated the following model.

$$unemp_t = 2.4 + 0.6 unemp_{t-1} - 0.2 unemp_{t-2} + \varepsilon_t.$$

where ε_t denotes an independent, identically-distributed process with $E(\varepsilon_t) = 0$ and $\text{Var}(\varepsilon_t) = \sigma^2$, constant. In our data set, $unemp_{T-1} = 5$ and $unemp_T = 4.5$. Compute the following forecast values.

a. Compute the forecast of $unemp_{T+1}$.

b. Compute the forecast of $unemp_{T+2}$.

c. Compute the limit of the forecast $unemp_{T+h}$ as h approaches infinity. [Hint: this is the unconditional mean of $unemp_t$.]

(6) [Forecasting, MA model: 9 pts] Suppose we have estimated the following model.

$$int_t = 4.5 + \varepsilon_t + 0.4 \varepsilon_{t-1}.$$

where ε_t denotes an independent, identically-distributed process with $E(\varepsilon_t) = 0$ and $\text{Var}(\varepsilon_t) = \sigma^2$, constant. In our data set, $\hat{\varepsilon}_T = -1.5$. Compute the following forecast values.

a. Compute the forecast of int_{T+1} .

b. Compute the forecast of int_{T+2} .

c. Compute the limit of the forecast int_{T+h} as h approaches infinity. [Hint: this is the unconditional mean of int_t .]

IV. CRITICAL THINKING: [4 pts] To investigate the possible effect of libraries on violent crime, the following regression was estimated using data on U.S. states¹:

$$y_i = \frac{-1091.}{(535)} + \frac{11.85}{(1.30)} x_i$$

where y_i denotes the number of violent crimes committed in state i in 2007 and x_i denotes the number of libraries in state i in the same year. Numbers in parentheses are standard errors of the coefficient estimates. The estimate for the slope coefficient is significantly greater than zero at 0.1 percent significance. So the number of violent crimes is clearly *positively correlated* with the number of libraries across states. Is this evidence that libraries *cause* violent crime? If yes, explain why. If no, explain why not and suggest a better way to estimate the regression equation using these state-level data.

[end of exam]

¹ These are actual least-squares estimates using data from the *Statistical Abstract of the United States*, 2010. $n=50$.