

U of T E-MAIL: _____@MAIL.UTORONTO.CA

FIRST (Given) NAME: _____

UNIVERSITY OF TORONTO

Faculty of Arts & Science

APRIL 2022 EXAMINATIONS

ECO220Y1Y: Introduction to Data Analysis and Applied Econometrics

Duration: 3 hours

Aids Allowed: A non-programmable calculator

Exam Reminders:

- Fill out your e-mail, name, and UTORid on the top of this page.
 - Do not begin writing the actual exam until the announcements have ended and the Exam Facilitator has started the exam.
 - As a student, you help create a fair and inclusive writing environment. If you possess an unauthorized aid during an exam, you may be charged with an academic offence.
 - Turn off and place all cell phones, smart watches, electronic devices, and unauthorized study materials in your bag under your desk. If it is left in your pocket, it may be an academic offence.
 - When you are done your exam, raise your hand for someone to come and collect your exam. Do not collect your bag and jacket before your exam is handed in.
 - If you are feeling ill and unable to finish your exam, please bring it to the attention of an Exam Facilitator so it can be recorded before leaving the exam hall.
 - In the event of a fire alarm, do not check your cell phone when escorted outside.

Special Instructions:

- Write your answers clearly, completely, and concisely. ***Your entire answer must fit in the designated space provided immediately after each question.*** No extra space/pages are possible. ***Write in PENCIL and use an ERASER as needed.*** Follow the answer guides, which end each question, and avoid excessively long answers.

Exam Format and Grading Scheme:

- This exam includes these **10 pages** plus the *Supplement*. There are **7 questions** (most with multiple parts) with varying point values worth a total of **120 points**.
 - The *Supplement* is 12 pages and contains graphs, tables, and other materials required for each exam question and the aid sheets. **After the exam begins, carefully DETACH the Supplement.** Remember, you must write your answers on the exam papers in the designated spaces: the Supplement will NOT be graded.

Students must hand in all examination materials at the end

(1) See the *Supplement for Question (1): Credit Line Utilization*.

(a) [5 pts] Using Table 1, roughly what fraction of participants have a credit line limit (TRY) within one standard deviation of the mean? Referencing relevant information in Table 1, explain your reasoning. Answer with 2 sentences.

(b) [4 pts] Given Figure 8 and the histogram of credit line utilization, would the standard deviation be larger for all customers (excluding participants) or for participants? How do you know? Answer with 2 sentences.

(c) [7 pts] Given Figure 8, how to best describe the distribution of available credit (TRY) for all customers (excluding participants)? Discuss the shape, units, and notable features in this context. Answer with 2 – 3 sentences.

(d) [10 pts] Given Table A.8, define events for credit line utilization for a randomly selected participant as:

- **Event A:** Utilization above 0% and below 25% in the one quarter prior to the experiment
- **Event B:** Utilization of at least 50% and less than 75% in the five quarters prior to the experiment

What are the numeric values of each of these five probabilities: $P(A)$, $P(B)$, $P(A \& B)$, $P(A | B)$, and $P(B | A)$? Answer with a quantitative analysis that shows your work.

(2) [4 pts] See the ***Supplement for Question (2): Working from Home***. For the number of days per week that people want to work from home, compute a good approximation of the mean. Answer with a quantitative analysis.

(3) See the ***Supplement for Question (3): Five Facts about Beliefs and Portfolios***.

(a) [2 pts] From Table 2, among the 46,419 people completing the survey, for the expected 1-year stock return the interquartile range is _____. Answer with a number and its units.

(b) [9 pts] The correlation between the expected 1-year stock return and the expected 10-year stock return is 0.304. Compute the 99% confidence interval estimate of the difference in expected returns for a 1-year versus 10-year horizon. Answer with a quantitative analysis.

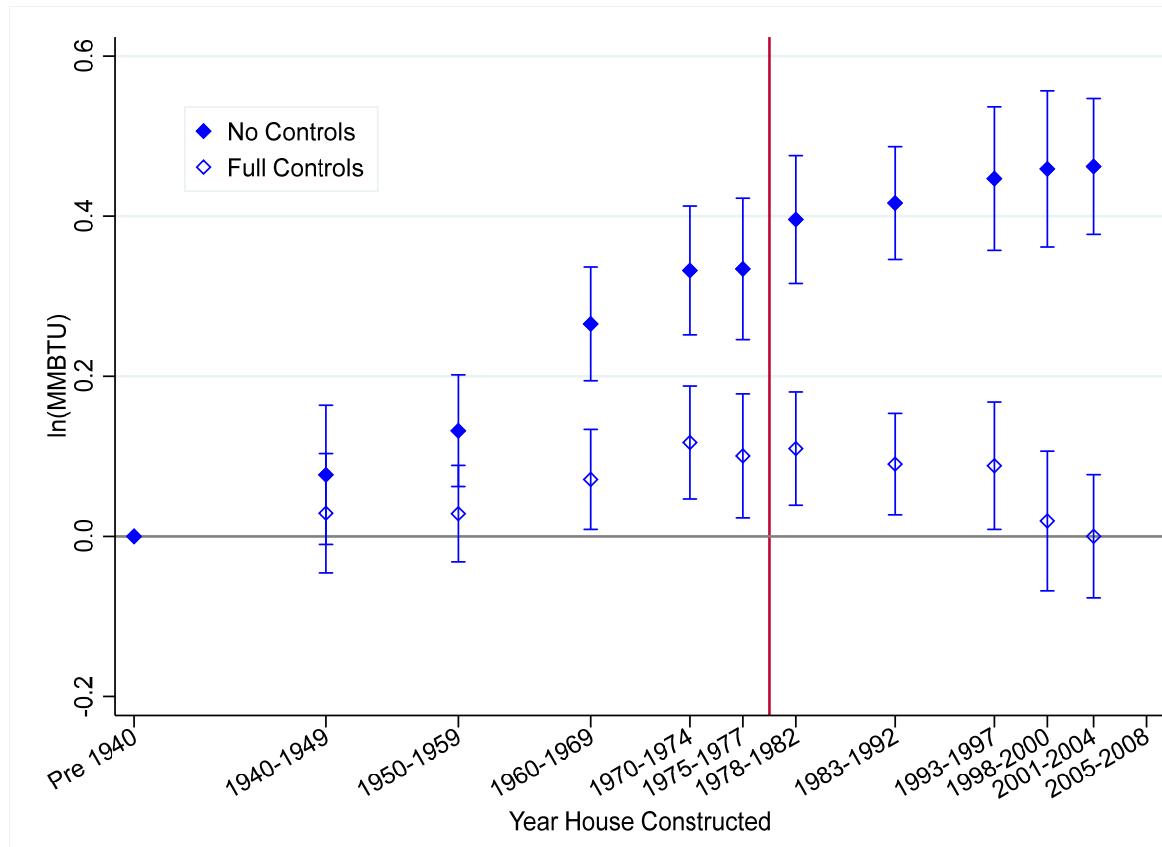
(4) [8 pts] See the ***Supplement for Question (4): Women's Downhill Skiing***. Why do the coefficient estimates in the row "start order" differ across Regressions #1, #2, and #3? Explain in this context. Answer with 3 – 4 sentences.

(5) See the ***Supplement for Question (5): California Energy.***

(a) [7 pts] In **Excel Output #1** find **2.97647E-70**. What is the hypothesis test? What does 2.97647E-70 imply? What does the null hypothesis mean *in this context?* Answer with hypotheses in formal notation & 2 – 3 sentences.

(b) [8 pts] Look at the “SS” column in **Excel Output #1** and **Excel Output #2**. *Why* is the total the same (1155.001745) in both? *Why* is the regression part smaller in the first regression: 89.904794 versus 375.9431875? Is this a small or big difference? Explain within this context. Answer with 3 – 4 sentences.

(c) [8 pts] Recall Figure 3 “Residential Electricity Use in California, Controlling for Characteristics.” The figure below is like Figure 3, but it uses the Excel regression output in the *Supplement*. It shows the 95% confidence interval estimates. Add the results for 2005-2008. Show your work. Answer with a quantitative analysis & complete the figure below.



(d) [12 pts] Is the coefficient on constr_05_08 in **Excel Output #2** statistically significant? What is the P-value? Next, *interpret* that coefficient. Answer with hypotheses in formal notation, a quantitative analysis & 2 sentences.

(6) See the *Supplement for Question (6): Partisan bias in inflation expectations, Part 1 of 2*.

(a) [11 pts] The authors claim: “The descriptive statistics in Table 3 show that the average inflation expectation is higher in red states (3.91%) than in blue states (3.63%), which is statistically significant at the 1% level.” Write the relevant hypotheses, graph the relevant *rejection region*, and compute the relevant test statistic to support their claim. Answer with hypotheses in formal notation, a fully labelled graph & a quantitative analysis.

(b) [10 pts] Is there a statistically significant difference between red and blue states in the fraction under 40 years old? Compute the P-value. Answer with hypotheses in formal notation, a quantitative analysis & 1 sentence.

(7) See the *Supplement for Question (7): Partisan bias in inflation expectations, Part 2 of 2*.

(a) [6 pts] In the excerpt, the researchers didn't say what is being controlled for when interpreting the results. In Column (2) of Table 4, correctly and completely interpret the coefficient estimate 0.390. Answer with 1 precise sentence.

(b) [4 pts] Using Column (1) of Table 4, what is the predicted inflation expectation in 2018 for a household in a red state in the Midwest, with a household head aged 48 years, with a high level of numeracy, who has no education beyond high school, and household income of \$44,000? Answer with a quantitative analysis that shows your work.

(c) [5 pts] Recall Table 3 from Question (6), which makes clear the reference (aka omitted) category for each suite of dummy variables. For Column (3) of Table 4, *interpret* the coefficient estimate for the row "Constant"? Answer with 1 precise sentence.

CAREFULLY DETACH THIS SUPPLEMENT FROM YOUR EXAM PAPERS NOW

This *Supplement* contains graphs, tables, and other materials required for each exam questions and the aid sheets (formulas and Normal, *t* and *F* statistical tables). Review all relevant materials for each question.

Supplement for Question (1): Consider Aydin (2022) "Consumption Response to Credit Expansions: Evidence from Experimental Assignment of 45,307 Credit Lines." See the excerpts, Table 1, Figure 8, and Table A.8.

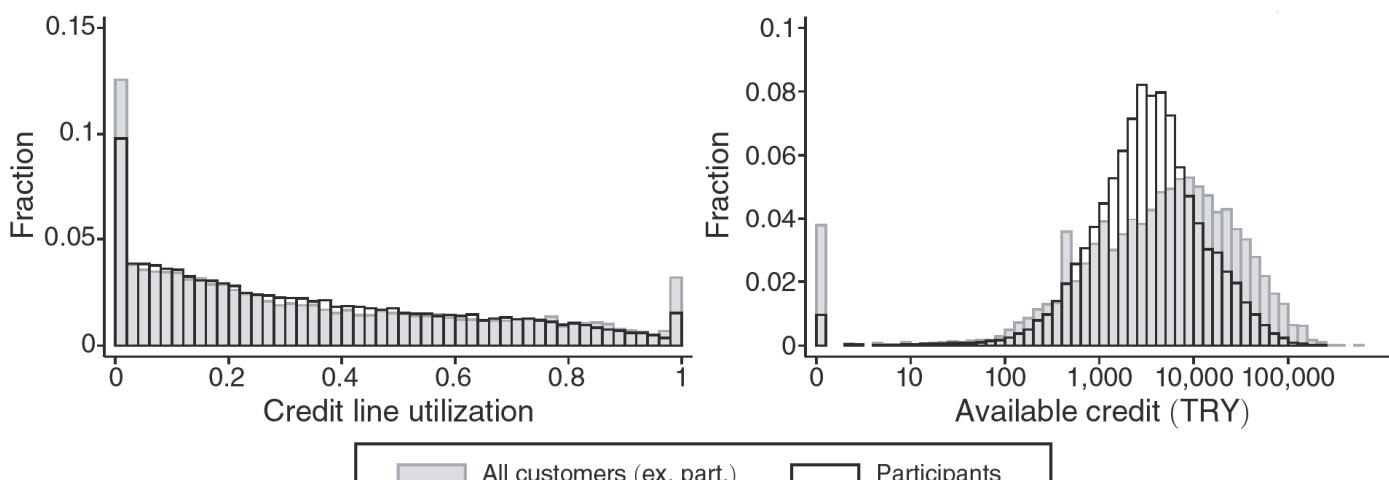
Excerpt, p. 7: How Participants Compare with the Typical Cardholder Table 1 displays summary statistics for the 45,307 participants in Panel A, and for a random sample of all credit line customers excluding participants ($N = 10,000$) in Panel B. Participants, compared with the universe of cardholders [which means the population of all cardholders], and on average, do not differ substantially in terms of age or labor income. However, the participants are not representative of the typical cardholder on several observable dimensions. For example, participants' median credit line across all banks is about 40 percent lower than that of the typical cardholder.

Table 1 – Summary Statistics

	Panel A. Participants					Panel B. Universe		
	N	Mean	SD	p10	p50	p90	Mean	p50
Age	45,307	37	10	26	35	50	41	40
Labor income (TRY)	17,690	2,465	2,423	943	1,600	5,111	2,292	1,426
<i>Credit lines (all banks)</i>								
Limit (TRY)	45,307	10,462	17,289	1,600	5,000	24,100	20,284	8,500
Debt (TRY)	45,307	3,446	8,619	94	1,277	6,978	6,220	1,983

Notes: Panel A is for $N = 45,307$ participants. It reports the mean, standard deviation, and 10th, 50th, and 90th percentiles. Panel B is for a random sample of the universe of all credit line customers excluding participants ($N = 10,000$). Statistics are from the quarter before the experiment: June 2014. The local currency of Turkish lira is abbreviated TRY. Limit (TRY) is the credit line limit – maximum total amount that can be borrowed – and Debt (TRY) is the amount that has been borrowed.

EXCERPT, p. 22: Figure 8 displays credit line utilization [the fraction of available credit that a person has already used] and unused limits in levels [how much more money in Turkish lira a person is still eligible to borrow] across all banks in the quarter before the experiment. The median and average credit line utilization are 0.27 and 0.34, respectively, with only one in ten of the participants utilizing more than 75 percent of their credit lines. The median and average available credit are 3,284 TRY and 7,016 TRY, corresponding to 1.33 and 2.85 times the average monthly post tax labor income.

**Figure 8:** Histograms of Credit Line Utilization and Available Credit (TRY)

Supplement for Question (1), continues the next page >>>

Supplement for Question (1), continued:

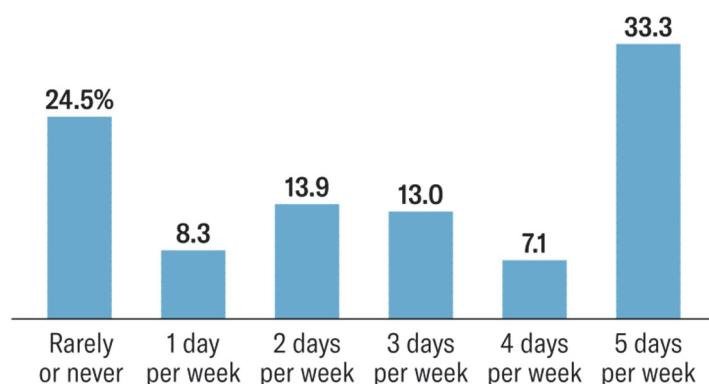
In Table A.8 below, the symbol “∅” is for households where credit line utilization is not in the data (i.e. is a missing value). For the ranges of values, curved parentheses “(” or “)” mean *not* including the endpoint and square brackets “[” or “]” mean *including* the endpoint. For example, [0.25, 0.50) is credit line utilization greater than or equal to 0.25 and strictly less than 0.50.

Table A.8: Utilization Transition Matrix:
Five quarters prior to the experiment to one quarter prior to the experiment

		One quarter prior				
Five quarters prior		=0	(0, 0.25)	[0.25, 0.50)	[0.50, 0.75)	[0.75, 1]
∅	319	1,165	1,071	879	653	
	0.08	0.29	0.26	0.22	0.16	
=0	764	1,434	570	341	191	
	0.23	0.43	0.17	0.10	0.06	
(0, 0.25)	960	9,370	3,393	1,433	551	
	0.06	0.60	0.22	0.09	0.04	
[0.25, 0.50)	345	4,040	3,248	1,957	954	
	0.03	0.38	0.31	0.19	0.09	
[0.50, 0.75)	166	1,758	1,959	1,787	1,115	
	0.02	0.26	0.29	0.26	0.16	
[0.75, 1]	107	899	1,186	1,370	1,322	
	0.02	0.18	0.24	0.28	0.27	
	2,661	18,666	11,427	7,767	4,786	
	0.06	0.41	0.25	0.17	0.11	

Notes: The 45,307 participants are allocated to $6 \times 5 = 30$ categories based on their utilization across all banks. Columns stand for utilization in the quarter before the experiment, the histogram of which is displayed in Figure 8. Rows stand for utilization four quarters prior to that. In each box, the first entry displays the number of participants, and the second row stands for the transition probabilities. Due to rounding, summed probabilities may not add up to 1.

Supplement for Question (2): In “Don’t Force People to Come Back to the Office Full Time” in the *Harvard Business Review* in August 2021, three economists discuss a survey of employees. The figure below summarizes the replies to one survey question: “After Covid, in 2022 and beyond, how often would you like to have paid workdays at home?”



This *Supplement* to the April 2022 ECO220Y1Y Final Exam will NOT be graded.

Supplement: Page 3 of 12

Supplement for Question (3): Consider Giglio et al. (2021) “Five Facts about Beliefs and Portfolios.” They use data from the GMSU-Vanguard survey that asks each investor about their expected 1-year stock return and their expected 10-year stock return, along with other survey questions. Below is part of Table 2 and an excerpt describing that part.

Excerpt, p. 1491: Table 2 shows summary statistics for the 46,419 survey responses. The average expected 1-year stock market return is 4.64 percent, while the average annualized expected 10-year stock market return is 6.64 percent. There is substantial heterogeneity in the expected 1-year stock market return across responses. At the tenth percentile of the distribution, individuals reported a 1-year expected stock return of -1 percent, while at the ninetieth percentile, they expected a return of 10 percent.

Table 2 – Summary Statistics: Survey Responses

	Mean	SD	P10	P25	P50	P75	P90
Expected 1Y stock return (percent)	4.64	6.08	-1	3	5	8	10
Expected 10Y stock return (annualized percent)	6.64	3.85	3	5	6	8	10

Notes: Table shows summary statistics of the answers for the 46,419 people completing the GMSU-Vanguard survey.

Supplement for Question (4): Recall women’s downhill skiing. Finish time is in seconds. Start order says when skiers race: 1 is first, 2 is second, and so on. FIS points measure skill as of Feb. 14, 2022, where 0 is best. In the 2022 Olympics, thirty-one athletes raced. Ester Ledecka lost control and is excluded as an outlier in the three regressions below.

Women’s Downhill Skiing Final: 2022 Beijing Olympics

	Dependent variable: Finish time (in seconds)		
	Regression #1	Regression #2	Regression #3
Start order	0.183 (0.034)	-0.343 (0.097)	-0.123 (0.104)
Start order squared		0.014 (0.003)	0.004 (0.004)
FIS points			0.109 (0.031)
Constant	91.579 (0.705)	94.905 (0.773)	93.335 (0.793)
Observations	30	30	30

Notes: Each column reports results from a separate OLS regression. Excludes Ester Ledecka. Standard errors are in parentheses.

Supplement for Question (5): Recall Levinson (2016) “How Much Energy Do Building Energy Codes Save? Evidence from California Houses.” The Residential Appliance Saturation Study (RASS) has many variables describing each house, its owners, the local climate, and the appliances. The key dependent variables are annual household electricity use in MMBTUs and annual household natural gas use in MMBTUs. Consider the RASS 2009 results for the ten counties in southern California: 4,438 houses. Consider the natural log of each household’s annual electricity use.

Dummy variables record when each house is constructed. For example, constr_75_77 is 1 if constructed from 1975-1977 and 0 otherwise. The reference (omitted) category is houses constructed before 1940. The variable cool_deg_days measures cooling degree days (CDDs) in 100s. The variable ln_sq_feet is the natural log of house size in square feet. The variable ln_num_res is the natural log of the number of residents. The variable central_ac is 1 if house has central air conditioning and is 0 otherwise. The variable central_ac_x_sq_ft is the product of central_ac and sq_feet divided by 1,000. There are two sets of Excel regression output on the next page. Some of the output has been intentionally erased.

Supplement for Question (5), continues the next page >>>

Supplement for Question (5), continued:

Excel Output #1: The dependent variable is the natural log of electricity use in MMBTUs in 2009 for 4,438 houses.

ANOVA					
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	11	89.904794	8.17316309	33.96349957	2.97647E-70
Residual	4426	1065.096951	0.240645493		
Total	4437	1155.001745			
	<i>Coefficients</i>	<i>Standard Error</i>			
Intercept	2.859580021	0.031025505			
constr_40_49	0.076898379	0.044331383			
constr_50_59	0.13193316	0.035576143			
constr_60_69	0.265411881	0.036268921			
constr_70_74	0.332168831	0.041025327			
constr_75_77	0.334065506	0.045026493			
constr_78_82	0.395822284	0.040719641			
constr_83_92	0.416355326	0.035928876			
constr_93_97	0.446899	0.04568472			
constr_98_00	0.458935587	0.049760181			
constr_01_04	0.462056557	0.043251337			
constr_05_08	0.425059892	0.046868555			

Excel Output #2: The dependent variable is the natural log of electricity use in MMBTUs in 2009 for 4,438 houses.

ANOVA					
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	16	375.9431875	23.49644922	133.3376047	0
Residual	4421	779.0585576	0.176217724		
Total	4437	1155.001745			
	<i>Coefficients</i>	<i>Standard Error</i>			
Intercept	2.35577861	0.031668966			
cool_deg_days	0.019928992	0.001862118			
ln_sq_feet	0.358512789	0.024221931			
ln_num_res	0.256123062	0.012132425			
central_ac	0.197658143	0.029827687			
central_ac_x_sq_ft	0.02294568	0.014315843			
constr_40_49	0.028935158	0.038002464			
constr_50_59	0.028344248	0.030740243			
constr_60_69	0.071133509	0.031877413			
constr_70_74	0.117263569	0.036032019			
constr_75_77	0.100545524	0.039472353			
constr_78_82	0.109658996	0.036096784			
constr_83_92	0.090265966	0.032361816			
constr_93_97	0.088252137	0.04057046			
constr_98_00	0.01921904	0.044510848			
constr_01_04	0.000146074	0.039238262			
constr_05_08	-0.095906238	0.042573558			

Supplement for Question (6): Bachmann et al. (2021) in “Partisan bias in inflation expectations” study how political beliefs affect peoples’ expectations about inflation. The New York Fed’s Survey of Consumer Expectations (SCE) asks a sample of heads of households what they think will happen to prices (i.e. inflation expectations). Each US state is either “red” (Republican leaning like Nebraska), “blue” (Democratic leaning like Massachusetts), or “swing” (sometimes goes Republican and sometimes goes Democratic like Michigan). Table 3 offers some descriptive statistics separately for red and blue states. Numeracy skills are the ability to work with numbers and data.

Table 3. Descriptive statistics

	Red	Blue
Inflation expectations, as a percent	3.91 (5.19)	3.63 (4.78)
Age		
Age under 40	0.2801	0.2690
Age 40–60	0.3930	0.4146
Age over 60	0.3269	0.3164
Numeracy skills		
Numeracy skills low	0.2954	0.2644
Numeracy skills high	0.7046	0.7356
Region		
West	0.0967	0.3305
Northeast	0.0003	0.4390
South	0.7250	0.0831
Midwest	0.1780	0.1474
Education		
High school degree or less	0.1284	0.1009
Some college	0.3600	0.2923
College	0.5116	0.6068
Household Income		
Income under 50 k	0.4082	0.2965
Income 50–100 k	0.3607	0.3356
Income over 100 k	0.2311	0.3679
Year		
2013	0.0896	0.1446
2014	0.1891	0.1990
2015	0.2069	0.1744
2016	0.1948	0.1854
2017	0.2111	0.2016
2018	0.1085	0.0950
Number of observations	18,103	23,629

Notes: Reports means. Standard deviation in parentheses for non-dummy variables.

The last row reports the number of observations for all variables, except for inflation expectation it is 17,931 for red states and 23,474 for blue states.

Supplement for Question (7): Continue with Bachmann et al. (2021) “Partisan bias in inflation expectations.” Table 4 offers some regression results. Also, there is an excerpt with the researchers’ discussion of Table 4.

Excerpt, p. 525: Table 4 shows the results of estimating the OLS model when we consider dummy variables for red and blue states (swing states are the reference category). When we consider the full sample (2013–2018) in Column (1), the blue-state dummy variable has a positive sign but does not turn out to be statistically significant. The red-state dummy variable has a positive sign and is statistically significant at the 1% level, indicating that inflation expectations in red states were around 0.28 percentage points higher than in swing states. In Column (2), we consider only the 2013–2016 period, when the Democrat Barack Obama was US president. The dummy variable for blue states has a negative sign and still lacks statistical significance; the point estimate of red states is still positive and statistically significant at the 1% level and is larger than in Column (1) when the full sample is used. Inflation expectations were around 0.39 percentage points higher in red states than in swing states. Compared to blue states, inflation expectations in red states are 0.46 percentage points higher, an effect that is statistically significant at the 1% level. In Column (3), we consider the 2017–2018 period only, when Republican Donald Trump was US president. The results change drastically: the dummy variable for blue states has a positive sign and is statistically significant at the 1% level, indicating that inflation expectations were around 0.27 percentage points higher in blue states than in swing states. By contrast, the dummy variable for red states lacks statistical significance.

Table 4. OLS regression results, with dummy variables for blue and red states

	Dependent variable: Inflation expectations, as a percent		
	Column (1): 2013 – 2018	Column (2): 2013 – 2016	Column (3): 2017 – 2018
Blue	0.041 (0.044)	-0.069 (0.052)	0.270*** (0.082)
Red	0.281*** (0.049)	0.390*** (0.059)	0.004 (0.084)
Age 40 to 60	0.653*** (0.042)	0.779*** (0.051)	0.353*** (0.074)
Age over 60	0.560*** (0.042)	0.746*** (0.051)	0.131* (0.076)
Numeracy High	-0.437*** (0.049)	-0.370*** (0.060)	-0.603*** (0.088)
Northeast	-0.453*** (0.054)	-0.532*** (0.066)	-0.278*** (0.095)
Midwest	-0.432*** (0.051)	-0.461*** (0.062)	-0.393*** (0.087)
South	-0.368*** (0.052)	-0.500*** (0.062)	-0.052 (0.093)
Some College	0.366*** (0.077)	0.082 (0.092)	1.026*** (0.139)
College	-0.169** (0.073)	-0.501*** (0.088)	0.598*** (0.131)
50 k < Income < 100 k	-0.380*** (0.045)	-0.391*** (0.054)	-0.366*** (0.080)
Income over 100 k	-0.792*** (0.045)	-0.798*** (0.055)	-0.785*** (0.079)
Year 2014	-0.154** (0.066)	-0.160** (0.066)	
Year 2015	-0.865*** (0.064)	-0.873*** (0.064)	
Year 2016	-0.929*** (0.064)	-0.932*** (0.064)	
Year 2017	-0.958*** (0.063)		
Year 2018	-0.781*** (0.072)		0.195*** (0.064)
Constant	4.793*** (0.106)	4.997*** (0.122)	3.406*** (0.165)
Observations	78,174	54,858	23,316

Notes: Standard errors in parentheses. ***Significant at the 1% level;

**Significant at the 5% level; *Significant at the 10% level.

This Supplement to the April 2022 ECO220Y1Y Final Exam will NOT be graded.

Supplement: Page 7 of 12

Sample mean: $\bar{X} = \frac{\sum_{i=1}^n x_i}{n}$ **Sample variance:** $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n-1} = \frac{\sum_{i=1}^n x_i^2}{n-1} - \frac{(\sum_{i=1}^n x_i)^2}{n(n-1)}$ **Sample s.d.:** $s = \sqrt{s^2}$

Sample coefficient of variation: $CV = \frac{s}{\bar{X}}$ **Sample covariance:** $s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{n-1} = \frac{\sum_{i=1}^n x_i y_i}{n-1} - \frac{(\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)}{n(n-1)}$

Sample interquartile range: $IQR = Q3 - Q1$ **Sample coefficient of correlation:** $r = \frac{s_{xy}}{s_x s_y} = \frac{\sum_{i=1}^n z_{x_i} z_{y_i}}{n-1}$

Addition rule: $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$ **Conditional probability:** $P(A|B) = \frac{P(A \text{ and } B)}{P(B)}$

Complement rules: $P(A^C) = P(A') = 1 - P(A)$ $P(A^C|B) = P(A'|B) = 1 - P(A|B)$

Multiplication rule: $P(A \text{ and } B) = P(A|B)P(B) = P(B|A)P(A)$

Expected value: $E[X] = \mu = \sum_{all x} x p(x)$ **Variance:** $V[X] = E[(X - \mu)^2] = \sigma^2 = \sum_{all x} (x - \mu)^2 p(x)$

Covariance: $COV[X, Y] = E[(X - \mu_X)(Y - \mu_Y)] = \sigma_{XY} = \sum_{all x} \sum_{all y} (x - \mu_X)(y - \mu_Y)p(x, y)$

Laws of expected value:

$$\begin{aligned} E[c] &= c \\ E[X + c] &= E[X] + c \\ E[cX] &= cE[X] \\ E[a + bX + cY] &= a + bE[X] + cE[Y] \end{aligned}$$

Laws of variance:

$$\begin{aligned} V[c] &= 0 \\ V[X + c] &= V[X] \\ V[cX] &= c^2 V[X] \\ V[a + bX + cY] &= b^2 V[X] + c^2 V[Y] + 2bc * COV[X, Y] \\ V[a + bX + cY] &= b^2 V[X] + c^2 V[Y] + 2bc * SD(X) * SD(Y) * \rho \end{aligned}$$

Laws of covariance:

$$\begin{aligned} COV[X, c] &= 0 \\ COV[a + bX, c + dY] &= bd * COV[X, Y] \\ where \rho &= CORRELATION[X, Y] \end{aligned}$$

Combinatorial formula: $C_x^n = \frac{n!}{x!(n-x)!}$ **Binomial probability:** $p(x) = \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x}$ for $x = 0, 1, 2, \dots, n$

If X is Binomial ($X \sim B(n, p)$) then $E[X] = np$ and $V[X] = np(1-p)$

If X is Uniform ($X \sim U[a, b]$) then $f(x) = \frac{1}{b-a}$ and $E[X] = \frac{a+b}{2}$ and $V[X] = \frac{(b-a)^2}{12}$

Sampling distribution of \bar{X} :

$$\begin{aligned} \mu_{\bar{X}} &= E[\bar{X}] = \mu \\ \sigma_{\bar{X}}^2 &= V[\bar{X}] = \frac{\sigma^2}{n} \\ \sigma_{\bar{X}} &= SD[\bar{X}] = \frac{\sigma}{\sqrt{n}} \end{aligned}$$

Sampling distribution of \hat{P} :

$$\begin{aligned} \mu_{\hat{P}} &= E[\hat{P}] = p \\ \sigma_{\hat{P}}^2 &= V[\hat{P}] = \frac{p(1-p)}{n} \\ \sigma_{\hat{P}} &= SD[\hat{P}] = \sqrt{\frac{p(1-p)}{n}} \end{aligned}$$

Sampling distribution of $(\hat{P}_2 - \hat{P}_1)$:

$$\begin{aligned} \mu_{\hat{P}_2 - \hat{P}_1} &= E[\hat{P}_2 - \hat{P}_1] = p_2 - p_1 \\ \sigma_{\hat{P}_2 - \hat{P}_1}^2 &= V[\hat{P}_2 - \hat{P}_1] = \frac{p_2(1-p_2)}{n_2} + \frac{p_1(1-p_1)}{n_1} \\ \sigma_{\hat{P}_2 - \hat{P}_1} &= SD[\hat{P}_2 - \hat{P}_1] = \sqrt{\frac{p_2(1-p_2)}{n_2} + \frac{p_1(1-p_1)}{n_1}} \end{aligned}$$

Sampling distribution of $(\bar{X}_1 - \bar{X}_2)$, independent samples:

$$\begin{aligned} \mu_{\bar{X}_1 - \bar{X}_2} &= E[\bar{X}_1 - \bar{X}_2] = \mu_1 - \mu_2 \\ \sigma_{\bar{X}_1 - \bar{X}_2}^2 &= V[\bar{X}_1 - \bar{X}_2] = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \\ \sigma_{\bar{X}_1 - \bar{X}_2} &= SD[\bar{X}_1 - \bar{X}_2] = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \end{aligned}$$

Sampling distribution of (\bar{X}_d) , paired ($d = X_1 - X_2$):

$$\begin{aligned} \mu_{\bar{X}_d} &= E[\bar{X}_d] = \mu_1 - \mu_2 \\ \sigma_{\bar{X}_d}^2 &= V[\bar{X}_d] = \frac{\sigma_d^2}{n} = \frac{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}{n} \\ \sigma_{\bar{X}_d} &= SD[\bar{X}_d] = \frac{\sigma_d}{\sqrt{n}} = \sqrt{\frac{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}{n}} \end{aligned}$$

Inference about a population proportion:

$$\text{z test statistic: } z = \frac{\hat{P} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \quad \text{CI estimator: } \hat{P} \pm z_{\alpha/2} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}$$

Inference about comparing two population proportions:

$$\text{z test statistic under Null hypothesis of no difference: } z = \frac{\hat{P}_2 - \hat{P}_1}{\sqrt{\frac{\hat{P}(1-\hat{P})}{n_1} + \frac{\hat{P}(1-\hat{P})}{n_2}}} \quad \text{Pooled proportion: } \bar{P} = \frac{x_1 + x_2}{n_1 + n_2}$$

$$\text{CI estimator: } (\hat{P}_2 - \hat{P}_1) \pm z_{\alpha/2} \sqrt{\frac{\hat{P}_2(1-\hat{P}_2)}{n_2} + \frac{\hat{P}_1(1-\hat{P}_1)}{n_1}}$$

Inference about the population mean:

$$\text{t test statistic: } t = \frac{\bar{X} - \mu_0}{s/\sqrt{n}} \quad \text{CI estimator: } \bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}} \quad \text{Degrees of freedom: } v = n - 1$$

Inference about a comparing two population means, independent samples, unequal variances:

$$\text{t test statistic: } t = \frac{(\bar{X}_1 - \bar{X}_2) - \Delta_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad \text{CI estimator: } (\bar{X}_1 - \bar{X}_2) \pm t_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$\text{Degrees of freedom: } v = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{1}{n_1-1}\left(\frac{s_1^2}{n_1}\right)^2 + \frac{1}{n_2-1}\left(\frac{s_2^2}{n_2}\right)^2}$$

Inference about a comparing two population means, independent samples, assuming equal variances:

$$\text{t test statistic: } t = \frac{(\bar{X}_1 - \bar{X}_2) - \Delta_0}{\sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}}} \quad \text{CI estimator: } (\bar{X}_1 - \bar{X}_2) \pm t_{\alpha/2} \sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}} \quad \text{Degrees of freedom: } v = n_1 + n_2 - 2$$

$$\text{Pooled variance: } s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}$$

Inference about a comparing two population means, paired data: (n is number of pairs and $d = X_1 - X_2$)

$$\text{t test statistic: } t = \frac{\bar{d} - \Delta_0}{s_d/\sqrt{n}} \quad \text{CI estimator: } \bar{X}_d \pm t_{\alpha/2} \frac{s_d}{\sqrt{n}} \quad \text{Degrees of freedom: } v = n - 1$$

SIMPLE REGRESSION:

$$\text{Model: } y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad \text{OLS line: } \hat{y}_i = b_0 + b_1 x_i \quad b_1 = \frac{s_{xy}}{s_x^2} = r \frac{s_y}{s_x} \quad b_0 = \bar{Y} - b_1 \bar{X}$$

$$\text{Coefficient of determination: } R^2 = (r)^2 \quad \text{Residuals: } e_i = y_i - \hat{y}_i$$

$$\text{Standard deviation of residuals: } s_e = \sqrt{\frac{SSE}{n-2}} = \sqrt{\frac{\sum_{i=1}^n (e_i - 0)^2}{n-2}} \quad \text{Standard error of slope: } s.e.(b_1) = s_{b_1} = \frac{s_e}{\sqrt{(n-1)s_x^2}}$$

Inference about the population slope:

t test statistic: $t = \frac{b_1 - \beta_{10}}{s.e.(b_1)}$ **CI estimator:** $b_1 \pm t_{\alpha/2} s.e.(b_1)$ **Degrees of freedom:** $v = n - 2$

Standard error of slope: $s.e.(b_1) = s_{b_1} = \frac{s_e}{\sqrt{(n-1)s_x^2}}$

Prediction interval for y at given value of x (x_g):

$$\hat{y}_{x_g} \pm t_{\alpha/2} s_e \sqrt{1 + \frac{1}{n} + \frac{(x_g - \bar{X})^2}{(n-1)s_x^2}} \quad \text{or} \quad \hat{y}_{x_g} \pm t_{\alpha/2} \sqrt{(s.e.(b_1))^2 (x_g - \bar{X})^2 + \frac{s_e^2}{n} + s_e^2}$$

Degrees of freedom: $v = n - 2$

Confidence interval for predicted mean at given value of x (x_g):

$$\hat{y}_{x_g} \pm t_{\alpha/2} s_e \sqrt{\frac{1}{n} + \frac{(x_g - \bar{X})^2}{(n-1)s_x^2}} \quad \text{or} \quad \hat{y}_{x_g} \pm t_{\alpha/2} \sqrt{(s.e.(b_1))^2 (x_g - \bar{X})^2 + \frac{s_e^2}{n}} \quad \text{Degrees of freedom: } v = n - 2$$

SIMPLE & MULTIPLE REGRESSION:

Model: $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \varepsilon_i$

$$SST = \sum_{i=1}^n (y_i - \bar{Y})^2 = SSR + SSE \quad SSR = \sum_{i=1}^n (\hat{y}_i - \bar{Y})^2 \quad SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$s_y^2 = \frac{SST}{n-1} \quad MSE = \frac{SSE}{n-k-1} \quad \text{Root MSE} = \sqrt{\frac{SSE}{n-k-1}} \quad MSR = \frac{SSR}{k}$$

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} \quad \text{Adj. } R^2 = 1 - \frac{SSE/(n-k-1)}{SST/(n-1)} = \left(R^2 - \frac{k}{n-1} \right) \left(\frac{n-1}{n-k-1} \right)$$

$$\text{Residuals: } e_i = y_i - \hat{y}_i \quad \text{Standard deviation of residuals: } s_e = \sqrt{\frac{SSE}{n-k-1}} = \sqrt{\frac{\sum_{i=1}^n (e_i - 0)^2}{n-k-1}}$$

Inference about the overall statistical significance of the regression model:

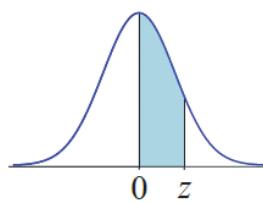
$$F = \frac{R^2/k}{(1-R^2)/(n-k-1)} = \frac{(SST-SSE)/k}{SSE/(n-k-1)} = \frac{SSR/k}{SSE/(n-k-1)} = \frac{MSR}{MSE}$$

Numerator degrees of freedom: $v_1 = k$ **Denominator degrees of freedom:** $v_2 = n - k - 1$

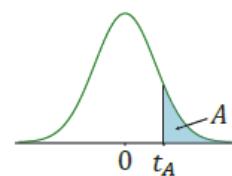
Inference about the population slope for explanatory variable j:

t test statistic: $t = \frac{b_j - \beta_{j0}}{s_{b_j}}$ **CI estimator:** $b_j \pm t_{\alpha/2} s_{b_j}$ **Degrees of freedom:** $v = n - k - 1$

Standard error of slope: $s.e.(b_j) = s_{b_j}$ (for multiple regression, must be obtained from technology)

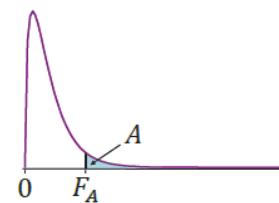


The Standard Normal Distribution:

Critical Values of Student t Distribution:

ν	$t_{0.10}$	$t_{0.05}$	$t_{0.025}$	$t_{0.01}$	$t_{0.005}$	$t_{0.001}$	$t_{0.0005}$	ν	$t_{0.10}$	$t_{0.05}$	$t_{0.025}$	$t_{0.01}$	$t_{0.005}$	$t_{0.001}$	$t_{0.0005}$
1	3.078	6.314	12.71	31.82	63.66	318.3	636.6	38	1.304	1.686	2.024	2.429	2.712	3.319	3.566
2	1.886	2.920	4.303	6.965	9.925	22.33	31.60	39	1.304	1.685	2.023	2.426	2.708	3.313	3.558
3	1.638	2.353	3.182	4.541	5.841	10.21	12.92	40	1.303	1.684	2.021	2.423	2.704	3.307	3.551
4	1.533	2.132	2.776	3.747	4.604	7.173	8.610	41	1.303	1.683	2.020	2.421	2.701	3.301	3.544
5	1.476	2.015	2.571	3.365	4.032	5.893	6.869	42	1.302	1.682	2.018	2.418	2.698	3.296	3.538
6	1.440	1.943	2.447	3.143	3.707	5.208	5.959	43	1.302	1.681	2.017	2.416	2.695	3.291	3.532
7	1.415	1.895	2.365	2.998	3.499	4.785	5.408	44	1.301	1.680	2.015	2.414	2.692	3.286	3.526
8	1.397	1.860	2.306	2.896	3.355	4.501	5.041	45	1.301	1.679	2.014	2.412	2.690	3.281	3.520
9	1.383	1.833	2.262	2.821	3.250	4.297	4.781	46	1.300	1.679	2.013	2.410	2.687	3.277	3.515
10	1.372	1.812	2.228	2.764	3.169	4.144	4.587	47	1.300	1.678	2.012	2.408	2.685	3.273	3.510
11	1.363	1.796	2.201	2.718	3.106	4.025	4.437	48	1.299	1.677	2.011	2.407	2.682	3.269	3.505
12	1.356	1.782	2.179	2.681	3.055	3.930	4.318	49	1.299	1.677	2.010	2.405	2.680	3.265	3.500
13	1.350	1.771	2.160	2.650	3.012	3.852	4.221	50	1.299	1.676	2.009	2.403	2.678	3.261	3.496
14	1.345	1.761	2.145	2.624	2.977	3.787	4.140	51	1.298	1.675	2.008	2.402	2.676	3.258	3.492
15	1.341	1.753	2.131	2.602	2.947	3.733	4.073	52	1.298	1.675	2.007	2.400	2.674	3.255	3.488
16	1.337	1.746	2.120	2.583	2.921	3.686	4.015	53	1.298	1.674	2.006	2.399	2.672	3.251	3.484
17	1.333	1.740	2.110	2.567	2.898	3.646	3.965	54	1.297	1.674	2.005	2.397	2.670	3.248	3.480
18	1.330	1.734	2.101	2.552	2.878	3.610	3.922	55	1.297	1.673	2.004	2.396	2.668	3.245	3.476
19	1.328	1.729	2.093	2.539	2.861	3.579	3.883	60	1.296	1.671	2.000	2.390	2.660	3.232	3.460
20	1.325	1.725	2.086	2.528	2.845	3.552	3.850	65	1.295	1.669	1.997	2.385	2.654	3.220	3.447
21	1.323	1.721	2.080	2.518	2.831	3.527	3.819	70	1.294	1.667	1.994	2.381	2.648	3.211	3.435
22	1.321	1.717	2.074	2.508	2.819	3.505	3.792	75	1.293	1.665	1.992	2.377	2.643	3.202	3.425
23	1.319	1.714	2.069	2.500	2.807	3.485	3.768	80	1.292	1.664	1.990	2.374	2.639	3.195	3.416
24	1.318	1.711	2.064	2.492	2.797	3.467	3.745	90	1.291	1.662	1.987	2.368	2.632	3.183	3.402
25	1.316	1.708	2.060	2.485	2.787	3.450	3.725	100	1.290	1.660	1.984	2.364	2.626	3.174	3.390
26	1.315	1.706	2.056	2.479	2.779	3.435	3.707	120	1.289	1.658	1.980	2.358	2.617	3.160	3.373
27	1.314	1.703	2.052	2.473	2.771	3.421	3.690	140	1.288	1.656	1.977	2.353	2.611	3.149	3.361
28	1.313	1.701	2.048	2.467	2.763	3.408	3.674	160	1.287	1.654	1.975	2.350	2.607	3.142	3.352
29	1.311	1.699	2.045	2.462	2.756	3.396	3.659	180	1.286	1.653	1.973	2.347	2.603	3.136	3.345
30	1.310	1.697	2.042	2.457	2.750	3.385	3.646	200	1.286	1.653	1.972	2.345	2.601	3.131	3.340
31	1.309	1.696	2.040	2.453	2.744	3.375	3.633	250	1.285	1.651	1.969	2.341	2.596	3.123	3.330
32	1.309	1.694	2.037	2.449	2.738	3.365	3.622	300	1.284	1.650	1.968	2.339	2.592	3.118	3.323
33	1.308	1.692	2.035	2.445	2.733	3.356	3.611	400	1.284	1.649	1.966	2.336	2.588	3.111	3.315
34	1.307	1.691	2.032	2.441	2.728	3.348	3.601	500	1.283	1.648	1.965	2.334	2.586	3.107	3.310
35	1.306	1.690	2.030	2.438	2.724	3.340	3.591	750	1.283	1.647	1.963	2.331	2.582	3.101	3.304
36	1.306	1.688	2.028	2.434	2.719	3.333	3.582	1000	1.282	1.646	1.962	2.330	2.581	3.098	3.300
37	1.305	1.687	2.026	2.431	2.715	3.326	3.574	∞	1.282	1.645	1.960	2.326	2.576	3.090	3.291

Degrees of freedom: ν

The *F* Distribution:

ν_1	1	2	3	4	5	6	7	8	9	10	11	12	15	20	30	∞
ν_2 Critical Values of <i>F</i> Distribution for $A = 0.10$:																
5	4.06	3.78	3.62	3.52	3.45	3.40	3.37	3.34	3.32	3.30	3.28	3.27	3.24	3.21	3.17	3.10
10	3.29	2.92	2.73	2.61	2.52	2.46	2.41	2.38	2.35	2.32	2.30	2.28	2.24	2.20	2.16	2.06
15	3.07	2.70	2.49	2.36	2.27	2.21	2.16	2.12	2.09	2.06	2.04	2.02	1.97	1.92	1.87	1.76
20	2.97	2.59	2.38	2.25	2.16	2.09	2.04	2.00	1.96	1.94	1.91	1.89	1.84	1.79	1.74	1.61
30	2.88	2.49	2.28	2.14	2.05	1.98	1.93	1.88	1.85	1.82	1.79	1.77	1.72	1.67	1.61	1.46
40	2.84	2.44	2.23	2.09	2.00	1.93	1.87	1.83	1.79	1.76	1.74	1.71	1.66	1.61	1.54	1.38
60	2.79	2.39	2.18	2.04	1.95	1.87	1.82	1.77	1.74	1.71	1.68	1.66	1.60	1.54	1.48	1.29
120	2.75	2.35	2.13	1.99	1.90	1.82	1.77	1.72	1.68	1.65	1.63	1.60	1.55	1.48	1.41	1.19
∞	2.71	2.30	2.08	1.94	1.85	1.77	1.72	1.67	1.63	1.60	1.57	1.55	1.49	1.42	1.34	1.00
ν_2 Critical Values of <i>F</i> Distribution for $A = 0.05$:																
5	6.61	5.79	5.41	5.19	5.05	4.95	4.88	4.82	4.77	4.74	4.70	4.68	4.62	4.56	4.50	4.36
10	4.96	4.10	3.71	3.48	3.33	3.22	3.14	3.07	3.02	2.98	2.94	2.91	2.85	2.77	2.70	2.54
15	4.54	3.68	3.29	3.06	2.90	2.79	2.71	2.64	2.59	2.54	2.51	2.48	2.40	2.33	2.25	2.07
20	4.35	3.49	3.10	2.87	2.71	2.60	2.51	2.45	2.39	2.35	2.31	2.28	2.20	2.12	2.04	1.84
30	4.17	3.32	2.92	2.69	2.53	2.42	2.33	2.27	2.21	2.16	2.13	2.09	2.01	1.93	1.84	1.62
40	4.08	3.23	2.84	2.61	2.45	2.34	2.25	2.18	2.12	2.08	2.04	2.00	1.92	1.84	1.74	1.51
60	4.00	3.15	2.76	2.53	2.37	2.25	2.17	2.10	2.04	1.99	1.95	1.92	1.84	1.75	1.65	1.39
120	3.92	3.07	2.68	2.45	2.29	2.18	2.09	2.02	1.96	1.91	1.87	1.83	1.75	1.66	1.55	1.25
∞	3.84	3.00	2.60	2.37	2.21	2.10	2.01	1.94	1.88	1.83	1.79	1.75	1.67	1.57	1.46	1.00
ν_2 Critical Values of <i>F</i> Distribution for $A = 0.01$:																
5	16.3	13.3	12.1	11.4	11.0	10.7	10.5	10.3	10.2	10.1	9.96	9.89	9.72	9.55	9.38	9.02
10	10.0	7.56	6.55	5.99	5.64	5.39	5.20	5.06	4.94	4.85	4.77	4.71	4.56	4.41	4.25	3.91
15	8.68	6.36	5.42	4.89	4.56	4.32	4.14	4.00	3.89	3.80	3.73	3.67	3.52	3.37	3.21	2.87
20	8.10	5.85	4.94	4.43	4.10	3.87	3.70	3.56	3.46	3.37	3.29	3.23	3.09	2.94	2.78	2.42
30	7.56	5.39	4.51	4.02	3.70	3.47	3.30	3.17	3.07	2.98	2.91	2.84	2.70	2.55	2.39	2.01
40	7.31	5.18	4.31	3.83	3.51	3.29	3.12	2.99	2.89	2.80	2.73	2.66	2.52	2.37	2.20	1.80
60	7.08	4.98	4.13	3.65	3.34	3.12	2.95	2.82	2.72	2.63	2.56	2.50	2.35	2.20	2.03	1.60
120	6.85	4.79	3.95	3.48	3.17	2.96	2.79	2.66	2.56	2.47	2.40	2.34	2.19	2.03	1.86	1.38
∞	6.63	4.61	3.78	3.32	3.02	2.80	2.64	2.51	2.41	2.32	2.25	2.18	2.04	1.88	1.70	1.00
ν_2 Critical Values of <i>F</i> Distribution for $A = 0.001$:																
5	47.2	37.1	33.2	31.1	29.8	28.8	28.2	27.6	27.2	26.9	26.6	26.4	25.9	25.4	24.9	23.8
10	21.0	14.9	12.6	11.3	10.5	9.93	9.52	9.20	8.96	8.75	8.59	8.45	8.13	7.80	7.47	6.76
15	16.6	11.3	9.34	8.25	7.57	7.09	6.74	6.47	6.26	6.08	5.94	5.81	5.54	5.25	4.95	4.31
20	14.8	9.95	8.10	7.10	6.46	6.02	5.69	5.44	5.24	5.08	4.94	4.82	4.56	4.29	4.00	3.38
30	13.3	8.77	7.05	6.12	5.53	5.12	4.82	4.58	4.39	4.24	4.11	4.00	3.75	3.49	3.22	2.59
40	12.6	8.25	6.59	5.70	5.13	4.73	4.44	4.21	4.02	3.87	3.75	3.64	3.40	3.14	2.87	2.23
60	12.0	7.77	6.17	5.31	4.76	4.37	4.09	3.86	3.69	3.54	3.42	3.32	3.08	2.83	2.55	1.89
120	11.4	7.32	5.78	4.95	4.42	4.04	3.77	3.55	3.38	3.24	3.12	3.02	2.78	2.53	2.26	1.54
∞	10.83	6.91	5.42	4.62	4.10	3.74	3.47	3.27	3.10	2.96	2.84	2.74	2.51	2.27	1.99	1.00

Numerator degrees of freedom: ν_1 ; Denominator degrees of freedom: ν_2