

Review of Regression

Joshua L Eubanks

University of Central Florida

Regression Equations

$$\hat{y} = \beta_0 + \beta_1 x_1 + \varepsilon$$

- ▶ \hat{y} is the dependent variable, also referred to as *response* or *predicted* variable
- ▶ x is the independent variable, also referred to as *explanatory* or *predictor* variable
- ▶ β_0 is the intercept
- ▶ β_1 is the slope
- ▶ ε is the random disturbance, which is unknown

Analysis of Regression Equations

Prediction and forecasting (Thing 1: Plug and Chug): plug in new observations into regression equation and calculate the predicted value of the dependent variable.

Marginal Analysis (Thing 2): Multiply the change in one independent variable by its slope to find the predicted change in the dependent variable.

Trend Analysis

When time is the independent variable the equation becomes:

$$\hat{y} = \beta_0 + \beta_1 \text{Time} + \varepsilon$$

- ▶ Time trends: upward, downward, or no trend
- ▶ Trend rates only apply to time trends
- ▶ Trend rate is the coefficient of the time variable, or the slope regression equation (β_1)

Simple vs Multiple Regression

- ▶ Simple Regression:

$$\hat{y} = \beta_0 + \beta_1 x_1 + \varepsilon$$

- ▶ Multiple Regression:

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$

- ▶ k = number of independent variables
- ▶ So $k = 1$ for Simple Regression
- ▶ $k > 1$ for Multiple Regression

Two Measures of Fit R^2 and SEE

- ▶ The regression is the best fitted equation, but how good is the best? To answer this we rely on two measures of fit:
 - ▶ $SEE = S$ = standard error (of the estimate), sometimes called the standard deviation of the regression
 - ▶ SEE is the Absolute Fit Measure that is in same units as the dependent variable
 - ▶ Thing 1 Prediction Margin of Error (MoE) is approximately $2 * SEE$ (for 95% cl)
 - ▶ $R^2 = R\text{-Squared} = R\text{-sqr}$
 - ▶ R^2 is the Relative measure of fit: $0\% \leq R^2 \leq 100\%$
 - ▶ Measure the percent of total variation in Y that is explained by the regression
 - ▶ Perfect Fit: $R^2 = 100\%$ all variation in Y is explained by the regression
 - ▶ No Fit: $R^2 = 0\%$ no variation in Y is explained by the regression

Adjusted R^2

- ▶ Comparing Different R^2 's can be like comparing Apples to Oranges
 - ▶ with more variables (larger k) you have an edge
 - ▶ with less data (smaller n) it affects the comparison
 - ▶ To avoid comparing apples to Oranges, use the same data but with different variables (e.g. one model with say 5 variables and another with 6)
- ▶ Adjusted R^2 corrects for differences in k and n
 - ▶ $\text{Adjusted } R^2 < R^2$: the gap is roughly proportional to k
 - ▶ The larger the n , the smaller is that gap
- ▶ Adjusted R^2 compares different regressions
 - ▶ Adjusted R^2 is often reported instead of R^2

Sample Regression v. Population Regression

- ▶ Sample regression: fitted equation from a sample

$$\hat{y} = b_0 + b_1x_1 + b_2x_2 + \dots + b_kx_k$$

- ▶ Population regression equation: unknown relationship of dependent variable y with the k independent variables

$$y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \dots + \beta_kx_k + \varepsilon$$

Where $\beta_1, \beta_2, \dots, \beta_k$ are population coefficients and ε is the random disturbance and the source of the unexplained variation in regression equation

- ▶ The expected value of the sample coefficient is the population coefficient: $\mathbb{E}(b_i) = \beta_i$

Summary of the Gauss Markov Theorem

Regression estimators are unbiased if the first 3 conditions are satisfied, and efficient if all five are satisfied:

- ▶ 1. parameters $\beta_0, \beta_1, \beta_2, \dots, \beta_k$ have constant values
 - ▶ Lesson: inefficient estimates if not a homogeneous populations or time periods
- ▶ 2. expected value of random disturbance ε is zero
 - ▶ Lesson: intercept term helps out when $\mathbb{E}(\varepsilon)$ is far from zero
- ▶ 3. Independent variables are not correlated with ε
 - ▶ Lesson: Watch out for Specification Bias dangers (see below)
- ▶ 4. standard deviation of ε , σ_ε , remains constant
 - ▶ Lesson: inefficient cross section estimates if data varies widely
- ▶ 5. Uncorrelated ε 's among different observations
 - ▶ Lesson: inefficient time series estimates if errors have a memory, which is referred to as autocorrelation (use Durbin Watson Statistic to test for autocorrelation)

Summary of the Gauss-Markov Theorem

Regression estimators are unbiased if the first 3 conditions are satisfied, and efficient if all five are satisfied. These are simple solutions to insure each condition is satisfied:

- ▶ 1. parameters have constant values
 - ▶ Solution: Study only homogeneous populations or time periods; use dummy variables to control for heterogeneity
- ▶ 2. expected value of random disturbance ε is zero
 - ▶ Solution: Always allow the software to calculate an intercept term (i.e. never force the slope coefficient to zero)
- ▶ 3. Independent variables are not correlated with ε
 - ▶ Solution: consider using ratios or percentages rather than raw data. Check for omitted variable bias, specification bias, multicollinearity

Summary Cont.

- ▶ 4. standard deviation of ε , σ_ε , remains constant
 - ▶ Solution: look for subgroups in data and analyze separately; use summary data (like the mean value) rather than the raw data
- ▶ 5. Uncorrelated ε 's among different observations
 - ▶ Solution: transform the data into differences or percent change; use a more sophisticated time series model like the autoregressive model which will model the error term's rate of decay

Model Specification and Specification Bias

- ▶ Modeling: only represents the most important factors in a model
- ▶ Use independent variables that are the best proxies for each of those important factors
- ▶ Modeling errors are from omitting less important factors
- ▶ Random disturbance ε is sum of these errors (plus the sum of the measurement errors of variables we do include)
- ▶ Omitting less important variables won't affect significance test results for important variables we keep in the model
- ▶ Omitting important variables that should be in the model reduces your fit (lower R^2 , higher SEE), and gives worse predictions
- ▶ But more importantly, you can no longer trust the t-test results if regression omits important independent variables that should be in model

Results of Specification Bias

- ▶ Biased coefficient estimates result from a misspecification of the regression model
- ▶ Specification bias may cause any of these problems:
 - ▶ Coefficient biased high or low by an unknown amount
 - ▶ Coefficient is the wrong sign (i.e., bias so large that it flips the sign of the coefficient)
 - ▶ Coefficient test gives wrong results: tests significant when it shouldn't, or doesn't test significant when it should

The Opposite of an Under-Specified model is an Over-Specified model

- ▶ Multicollinearity: some independent variables contain some of the same information
- ▶ Having some of the same information results in correlation between two or more independent variables
- ▶ Common amongst related proxy variables
- ▶ Common in time series analysis as variables tend to move together

Extreme Case of Multicollinearity and Violation of Gauss-Markov Assumption 3

- ▶ Can express one independent variable as linear combination of rest: amount of information < number of indep variables
- ▶ e.g. Height in inches and feet ($HtInches = 12 HtFt$) - $Wt = -380 + 96 HtFt$ or $Wt = -380 + 8 HtInch$ - $HtFt = 5$, $Wt = 100$ and $HtInch = 60$, $Wt = 100$
- ▶ Perfect collinearity: include both by allocating coefficients
 - ▶ $Wt = -380 + 8 HtInch + 0 HtFt$
 - ▶ $Wt = -380 + 0 HtInch + 96 HtFt$
 - ▶ $Wt = -380 + 4 HtInch + 48 HtFt$
 - ▶ $Wt = -380 + 16 HtInch - 96 HtFt$
 - ▶ $Wt = -380 + 88 HtInch - 960 HtFt$
- ▶ Infinite number of equivalent versions, choice arbitrary

Perfect Multicollinearity

- ▶ Inflated coefficients (average of $\pm\infty$)
- ▶ Explanation: too many variables for amount of information in the data
- ▶ Result: Arbitrary regression coefficients, least squares explodes (computational singularity)
- ▶ Regression software saves us, fortunately
 - ▶ deletes last independent variable in collinear list
 - ▶ Equivalent of “Hey, dummy, you’re repeating yourself.”
- ▶ Treatment: avoid or remove redundant information from your model, e.g. $\text{NetWorth} = \text{Assets} - \text{Debt}$

Garden Variety Multicollinearity

- ▶ When two (or more) independent variable(s) is nearly a linear combination of the rest
- ▶ Simplest situation is if two independent variables are highly correlated
- ▶ In effect, the two independent variables provide approximately the same information
- ▶ More complex, general case is a subset of 3+ independent variables almost linearly related
- ▶ Warning: pairwise correlations may not be large when more than 2 variables involved

Problems with Severe Multicollinearity

- ▶ Least squares can't disentangle which regressors explain what
- ▶ Potential Damage of Multicollinearity:
 - ▶ inflated standard errors of regression coefficients, reduced t-ratios, high p-values
 - ▶ regression coefficients: arbitrary, volatile, unreliable
 - ▶ important regressors may not test significant

Can't Blame The Following Problems on Multicollinearity

- ▶ Entire model doesn't test significant (F test)
- ▶ Model has a poor fit (low R^2)
- ▶ Inaccurate and imprecise predictions or forecasts from the fitted regression equation
- ▶ Failure to test significant of regressors that are not part of the subset of collinear variables

Alternative Explanations that share similar symptoms

- ▶ Obvious alternative if not statistical significance: those regressors really do not matter
- ▶ Unjustified to blame multicollinearity for nonsignificance if correlations and VIFs not high
- ▶ VIF stands for variance inflation factor, which quantifies the severity of multicollinearity in an ordinary least squares regression analysis. $VIF=1$ no multicollinearity; $1 < VIF < 5$ multicollinearity not a concern; $VIF > 5$ multicollinearity “watch”, $VIF > 10$ multicollinearity “warning”
- ▶ VIF is sensitive to the number of independent variables. With small number of independent variables a $VIF = 5$ indicates a potential problem, whereas a $VIF = 5$ with a large number of independent variables does not.

Multicollinearity is not a problem unless it is a problem

- ▶ Not damaging if independent variables test significant despite high correlation
 - ▶ Just because correlations are high does not necessarily mean there is damage to test results
- ▶ Not damaging if primary goal is Thing 1 prediction or forecasting, multicollinearity not an issue
- ▶ Not damaging if collinearity is confined to “control variables” you’re not planning to test individually anyway

How to Treat Multicollinearity

- ▶ Avoid poorly defined variables and poorly designed model specification:
 - ▶ Avoid using raw data: use ratios or percents instead of totals and subtotals
 - ▶ Ex.: include population and %Elderly, not population and #Elderly
 - ▶ Make default a common category to minimize correlations
 - ▶ Ex.: use majority population as your default, so in Orlando white is the default and %Latinos, %Black, and %Asian are included in the model

What not to do about Multicollinearity

Do not remove legitimately specified regressors

- ▶ Omitting important regressors risks specification bias and reduces overall fit
 - ▶ “Out of frying pan, into the fire”
 - ▶ Resulting bias makes test results questionable or worthless

Fear of Multicollinearity Led to Misspecified Model that Omits Number of Times Ranked (Top25)

Correlation Table of independent variables

	Win%10Yr	Enrollmt
Enrollmt	0.280	
Top25	0.663	0.287

$$\text{AttendAve} = -60 + 1.5 \text{ Win%10Yr} + 0.64 \text{ Enrollmt}$$

Predictor	Coef	SE Coef	T	P
Constant	-59.75	12.26	-4.88	0.000
Win%10Yr	1.5245	0.2247	6.78	0.000 ← signif
Enrollmt	0.6438	0.1625	3.96	0.000

$$S = 16.1260 \quad R-Sq = 48.4\% \quad R-Sq(\text{adj}) = 47.2\%$$

Significance of Win Pct. Vanishes and Fit Way Up If Top25 Added in Correct Specification

The regression equation is

$$\text{AttendAv} = 0.80 + 0.29 \text{ Win%10Yr} + 0.49 \text{ Enrollmt} \\ + 4.7 \text{ Top25CNN}$$

Predictor	Coef	SE Coef	T	P
Constant	0.80	10.71	0.07	0.941
Win%10Yr	0.2874	0.2046	1.40	0.164 ← not
Enrollmt	0.4876	0.1157	4.22	0.000
Top25CNN	4.7376	0.4960	9.55	0.000

$$S = 11.3634 \quad R-\text{Sq} = 74.6\% \quad R-\text{Sq}(\text{adj}) = 73.8\%$$

Bad Gas Price Model Omits Dow Index (DJIA): fear high Correlation with Other Variables

Correlations among independent variables

	P CRUDE	OPER%	PPI
OPER%	-0.501		
PPI	0.076	0.595	
DJIA	-0.396	0.765	0.700

$$P \text{ GAS} = -2.0 + 2.2 P \text{ CRUDE} + 0.26 \text{ OPER\%} + 0.34 \text{ PPI}$$

Predictor	Coef	SE Coef	T	P
Constant	-1.97	14.74	-0.13	0.894
P CRUDE	2.2414	0.1700	13.19	0.000
OPER%	0.2596	0.2305	1.13	0.264<-not signif
PPI	0.3395	0.1611	2.11	0.039

$$S = 4.65032 \quad R-Sq = 82.5\% \quad R-Sq(\text{adj}) = 81.8\%$$

Dow Improves Fit, Avoids Specification Bias: Oper% and PPI proxy for Dow's Inverse Effect on Gas Prices

The regression equation is

$$\begin{aligned}P \text{ GAS} = & - 46 + 2.0 P \text{ CRUDE} + 0.53 \text{ OPER\%} \\& + 0.74 \text{ PPI} - 0.0077 \text{ DJIA}\end{aligned}$$

Predictor	Coef	SE Coef	T	P
Constant	-46.14	17.58	-2.62	0.011
P CRUDE	1.9994	0.1668	11.99	0.000
OPER%	0.5268	0.2207	2.39	0.020
PPI	0.7413	0.1794	4.13	0.000
DJIA	-0.007715	0.001985	-3.89	0.000

$$S = 4.23 \quad R-Sq = 85.7\% \quad R-Sq(\text{adj}) = 84.9\%$$

Wait on Mortgage Loan Approval: Bad Fit, Low Correlations, cannot blame Multicollinearity

Wait on Mortgage Loan Approval: Bad Fit, Low Correlations, cannot blame Multicollinearity

Correlations among independent variables **NONE HIGH!**

	INT	INC	LOAN
INC	-0.116		
LOAN	-0.235	0.347	
LATE30D	-0.035	0.299	0.054

$$\text{WAIT} = 7.8 - 0.36 \text{ INT} + 0.018 \text{ INC} - 0.026 \text{ LOAN} + 0.19 \text{ LATE30D}$$

Predictor	Coef	SE Coef	T	P
Constant	7.838	3.140	2.50	0.016
INT	-0.3571	0.2986	-1.20	0.237
INC	0.01816	0.02399	0.76	0.453
LOAN	-0.02630	0.01856	-1.42	0.162
LATE30D	0.1898	0.1053	1.80	0.077

$$S = 3.54912 \quad R-Sq = 12.3\% \quad R-Sq(\text{adj}) = 5.7\%$$

Summary of Practical Model Design: Basic Principles and Case Applications

- ▶ Include intercept so fit not forced thru origin: (Assump. 2)
- ▶ Don't mix apples with oranges in same model (Assump. 1)
- ▶ Avoid specification bias, wrong/suspect results (Assump. 3)
 - ▶ Include all "important" independent variables even if not interested in their effect or you fear multicollinearity damage
- ▶ Multicollinearity, at worst, damages only significance tests
 - ▶ For low correlations among indep. vars., no damage to tests
 - ▶ Causes no damage to forecasts, fit, or to uncorrelated variables
 - ▶ Only damaging if damaging: if variables don't test significant

Bad Modeling Practices

- ▶ Discard Variables Simply Because They Aren't Significant
 - ▶ Not significant: means $p \geq \alpha$, but variable may still matter
 - ▶ Ex., $p = 0.30$ still is a 70% chance the variable is related
 - ▶ Variables are in model because they belong to control for their effects
 - ▶ May not test significant due to multicollinearity
- ▶ Torturing Data: Try All Possible Models and Number of Observations Until You Get the "Best" Model (i.e. p-hacking)
 - ▶ "Data Dredging" to find maximum fit often called unethical!
 - ▶ Data mining finds sample patterns of individuals that may not be present in population
 - ▶ Variables are in the model because they belong there to control for effects
 - ▶ Clients won't believe results if their favorite variable isn't in the model

Functional Forms

- ▶ Misspecification may occur even if you don't omit important variables
- ▶ Another type of misspecification is assuming a linear function when it's actually nonlinear

Tricking Linear regression into running nonlinear functions

Principle: CLS model (Classical Least Square) is linear in coefficients, not in variables

- ▶ Create (transform) data columns as new variables from other columns -e.g. quadratic, reciprocal, square root, log or ln (natural logs), product/interactive, ratio
- ▶ multiplicative production & demand functions transformed by natural logs
- ▶ If can't transform nonlinear to least squares, must use maximum likelihood estimation instead

Functional Forms and Marginal Change

Table 7.1 Summary of Alternative Functional Forms

Functional Form	Equation (one X only)	The Change in Y when X Changes
Linear	$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$	If X increases by one unit, Y will change by β_1 units.
Double-log	$\ln Y_i = \beta_0 + \beta_1 \ln X_i + \epsilon_i$	If X increases by one percent, Y will change by β_1 percent. (Thus β_1 is the elasticity of Y with respect to X.)
Semilog ($\ln X$)	$Y_i = \beta_0 + \beta_1 \ln X_i + \epsilon_i$	If X increases by one percent, Y will change by $\beta_1/100$ units.
Semilog ($\ln Y$)	$\ln Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$	If X increases by one unit, Y will change by roughly $100\beta_1$ percent.
Polynomial	$Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \epsilon_i$	If X increases by one unit, Y will change by $(\beta_1 + 2\beta_2 X)$ units.

Consider Non-Linear Function if Justified by Theory, Experience, Logic, Previous Studies

- 1: Using Log of a variable - E.g. Population, enrollment, sales vary too much - E.g. Income directly related to education, but linear not logical (impact less at higher incomes)

$$Educ = \beta_0 + \beta_1 FamInc + \varepsilon$$

- So regress with a log-linear functional form

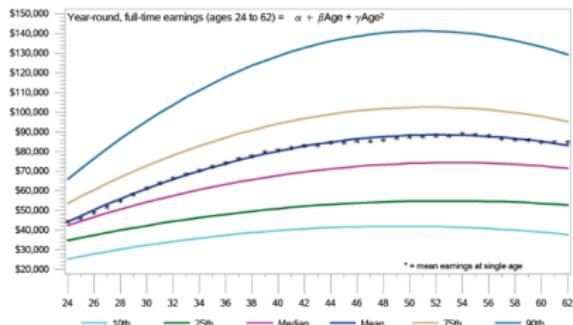
$$Educ = \beta_0 + \beta_1 Log10FamInc + \varepsilon$$

- Where define $Log10FamInc = \log_{10}(FamInc)$

Specify Quadratic if Expecting Maximum or Minimum

- ▶ Independent variable with too much of good thing → maximum
- ▶ E.g. Extra *Adv* helps less as the advertising budget increases (i.e. the law of diminishing returns):
 - ▶ $\text{NetInc} = \beta_0 + \beta_1 \text{Adv} + \beta_2 \text{AdvSq} + \varepsilon$
 - ▶ Where AdvSq is a transformed variable: $\text{AdvSq} = \text{Adv}^2$
- ▶ If both *Adv* and *AdvSq* test significant, then *NetInc* is a maximum at $\text{Adv} = -\frac{b_1}{2b_2}$
- ▶ Marginal change (thing 2) for $\Delta \text{Adv} = 1$ is:
$$\Delta \text{NetInc} = (b_1 + 2b_2 \text{Adv})$$
 - ▶ So, if $b_1 = 40$, $b_2 = -0.25$,
 - ▶ then max *NetInc* at $\text{Adv} = (-40 / -0.5) = 80$
 - ▶ so at max $\Delta \text{NetInc} = (40 + 2 * -.25 * 80) = 0$
- ▶ What if your *Adv* budget is 60 and you want to ask for more?
 - ▶ so at each additional dollar to *Adv* buget causes *NetInc* to change by
 - ▶ $\Delta \text{NetInc} = (40 + 2 * -.25 * 60) = 10$

Case Study Nonlinear Specified as Quadratic: Age Earnings Profile, Master degree



Conforming age-earnings regression results										
Level	α	β	γ	State adjustments at mean			State adjustments at mean			
				US	Pop	ACS Resp	10th	25th	Median	Mean
10th	-67,958	5,015,47	-52,449	9,865	12,936	1,423	-26.1	40.6	-36.9	FL 85.0
25th	-30,566	6,360,33	-44,092	2,511	123,94	11,451	-33.3	51.3	-43.9	GA 94.3
Median	-111,514	8,623,02	-44,559	2,711	134,04	13,573	-41.1	64.3	-53.8	HI 94.3
Mean	-205,739	13,401,87	-129,470	4,376	207,60	23,371	-47.0	64.6	-54.6	AZ 89.5
75th	+811,429	12,900,43	+123,312	3,001	180,24	22,223	+48.0	69.0	+93.3	CA 120.8
90th	-286,486	18,711,163	-11,494	61,195	8,150	-24.9	29.6	-19.9	CO 100.0	NJ 120.8
Cross-tabulation of earnings by age group										
Age	US	Pop	ACS Resp	10th	25th	Median	Mean	75th	90th	FL 85.0
18 & over	3,561,045	182,583	42,479	62,694	93,239	118,808	134,678	205,851	109,00	OH 90.0
18 to 62	3,412,525	173,973	42,853	62,656	93,239	118,451	134,678	203,813	109,00	AL 82.5
18 to 24	18,788	896	20,379	30,878	50,559	53,548	69,028	87,487	91.7	ND 87.2
18 to 25	18,788	896	20,379	30,878	50,559	53,548	69,028	87,487	91.7	AK 83.6
20 to 24	18,788	896	20,379	30,878	50,559	53,548	69,028	87,487	91.7	AR 80.5
20 to 25	18,788	896	20,379	30,878	50,559	53,548	69,028	87,487	91.7	AZ 87.1
25 to 34	797,338	37,588	86,446	107,39	130,378	154,559	153,548	160,028	177,99	CA 124.0
25 to 29	290,528	13,316	29,204	42,470	60,671	66,589	62,341	104,520	110,00	NH 100.0
30 to 34	300,810	24,270	45,141	33,871	77,088	89,398	104,320	143,007	152,00	NJ 118.8
35 to 39	1,000,000	1,000,000	1,000,000	1,000,000	1,000,000	1,000,000	1,000,000	1,000,000	1,000,000	DC 112.6
35 to 39	548,751	26,672	46,317	62,712	91,007	110,358	125,423	180,120	194,75	TX 98.4
40 to 44	541,897	27,018	49,913	70,059	102,429	128,467	140,327	212,378	173,00	GA 90.4
45 to 54	967,099	51,136	51,463	74,107	107,043	140,679	166,779	252,797	211,19	OR 87.5
45 to 54	538,047	44,411	50,975	54,209	79,213	99,294	102,213	122,297	114,27	VA 100.5
55 to 64	448,013	24,411	50,975	54,209	79,213	99,294	102,213	122,296	125,997	ID 104.7
55 to 64	538,661	30,701	45,599	68,029	103,598	134,097	153,586	242,485	195	DE 100.0
55 to 59	369,224	20,988	47,789	71,152	104,864	137,530	151,398	252,797	205,00	NY 111.4
55 to 59	68,462	8,721	42,479	47,602	99,097	126,617	145,038	217,406	181,00	DC 125.2
60 to 62	169,437	14,808	41,808	63,757	99,294	120,705	145,038	215,580	164,00	NY 89.7
63 to 67	134,345	8,711	41,808	63,757	99,294	120,705	145,038	215,580	164,00	WI 95.4
68 to 74	12,264	748	40,448	103,599	140,310	154,712	154,712	296,485	193,00	MS 71.8
75 & over	n/a	51	10,870	51,463	77,194	124,742	156,734	254,873	139,00	WV 81.9

Notes:

All national estimates are produced from ACS non-allocated populations of 28,000 or more wage & salary employees who do not have income from self-employment, retirement, Social Security, Supplemental Security, or public assistance. The indicator 'n/r' signals that the estimates for the subgroup are not reliable. The indicator 'n/d' signals that no data is available within the subgroup. Care should be exercised in using and interpreting earnings values for small subgroups of the population.

Notes:

All national estimates are produced from ACS non-allocated populations of 28,000 or more wage & salary employees who do not have income from self-employment, retirement, Social Security, Supplemental Security, or public assistance. The indicator 'n/r' signals that the estimates for the subgroup are not reliable. The indicator 'n/d' signals that no data is available within the subgroup. Care should be exercised in using and interpreting earnings values for small subgroups of the population.

Testing the Fit of a Regression Model

$$FTE = \alpha + \beta \text{Age} + \gamma \text{Age}^2 \quad \text{where FTE is full time earnings}$$

Why would you use a quadratic functional form? Just look at the plots on previous page.

Why does the study divide individuals by Gender, Educational Attainment and by percentiles of FTE? (Gauss Markov #4. standard deviation of ϵ , σ_ϵ , remains constant)

For Men in 90 percentile (in top 10%):

$$FTE = \alpha + \beta \text{Age} + \gamma \text{Age}^2 = -286,486 + 18,712(\text{Age}) - 162(\text{Age}^2)$$

(29.6) (-20) (t-ratios)

For Women in 90 percentile (in top 10%):

$$FTE = \alpha + \beta \text{Age} + \gamma \text{Age}^2 = -126,559 + 10,471(\text{Age}) - 102(\text{Age}^2)$$

(50) (-40) (t-ratios)

Do both β and γ test significant at 95% cl?

Yes, for both models: $|t-ratio| \geq 2$ for both Age and Age^2

So for Men- ETE peaks at what age? Max ETE where Age = $-B/2A = -18.712/2(-162) = 57.75$

Max FTE where Age = - β /2 γ = -10.471/2(-102) = 51.33

For a 52 year old women, each additional year of age adds how much to FTE?

$$\Delta FTE = (\beta + 2\gamma Age) = 10.471 + (2 * -102 * 52) = -137$$

For a 52 year old man, each additional year of age adds how much to FTE?

$$\Delta FTE = (\beta + 2\gamma Age) = 18.712 + (2 * -162 * 52) = 1.864$$

When to consider Nonlinear, if regressors are expected to interact

- ▶ product of 2 variables, e.g., slope dummy (see below)
- ▶ Synergism (positive interaction): more than sum of parts
- ▶ Negative interaction: conflict, too many cooks spoil the broth
- ▶ Slope dummy: product of quantitative and dummy var
 - ▶ E.g. recall men earnings: $\text{Wage} = \beta_0 + \beta_1 \text{Exp} + \beta_2 \text{Male} + \varepsilon$
 - ▶ vs Male raises higher: $\text{Wage} = \beta_0 + \beta_1 \text{Exp} + \beta_2 \text{ExpMale} + \varepsilon$
 - ▶ use slope dummy: $\text{ExpMale} = \text{Exp}$ if man, = 0 if woman
 - ▶ Same as 2 equations: $\text{Wage} = \beta_0 + \beta_1 \text{Exp} + \varepsilon$ for women
 - ▶ and $\text{Wage} = \beta_0 + (\beta_1 + \beta_2) \text{Exp} + \varepsilon$ for men
 - ▶ Interactive because of the product: $\text{ExpMale} = \text{Exp} * \text{Male}$
 - ▶ Slope dummy: slope is β_1 for women, $(\beta_1 + \beta_2)$ for men

When to consider Nonlinear, if regressors are expected to interact (continued)

- ▶ Dummy variable interaction Ex:

$$\text{Salary} = \beta_0 + \beta_1 \text{MBA} + \beta_2 \text{Eng} + \beta_3 \text{Both} + \varepsilon$$

- ▶ Where MBA and Engr are college degree dummies and
 $\text{Both} = \text{MBA} * \text{Eng}$ which equals 1 if have both degrees
 - ▶ E.g.: $\text{Salary} = 30 + 25\text{MBA} + 15\text{Eng} + 20\text{Both} + \varepsilon$

Means that MBAs earn \$25K more, engineers \$15K more, but with both degrees, earn $(25 + 15 + 20) = \$60K$ more But if $b_3 < 0$, negative interaction, whole < sum of parts

Multiplicative Functional Form

- ▶ $Y = aX^bZ^c$ nonlinear, so is linear regression hopeless? Not hopeless as you can easily transform to linear by taking natural logs. Recall the rules for logs:

$$\ln(xz) = \ln x + \ln z \text{ and } \ln(x^a) = a \ln x$$

- ▶ So transform $Y = aX^bZ^c$ into $\ln Y = \ln a + b \ln X + c \ln Z$
 - ▶ Note: can't compare R²'s with linear because different dep var!
- ▶ For demand, regress logarithms so that coefficients are elasticities

$$LQCar = 3.0 - 2LPCar - 0.5LPGas + 0.4LPlmprt + 4LInc$$

- ▶ Price elasticity $EP = -2$, +5% price cuts demand 2(5%) = 10%
- ▶ Elasticity for gas price, import price, income = -0.5, +0.4, +4

Other Functional Forms

- ▶ Ratios: e.g., financial ratios, percent Latino
- ▶ %Change or Difference: especially time series
- ▶ Reciprocal or square root: alternative to logs

When to Consider Nonlinear

- ▶ Indicated by theory or logic (e.g., log, quadratic)
- ▶ Variable collected over wide enough range of values
- ▶ Enough degrees of freedom: n much larger than k
- ▶ Easier to interpret results (e.g., demand elasticities)
- ▶ error analysis (residual plot) or variable plot