MULTIPLE REGRESSION

Chapter 4

Topics

- 1. What is regression analysis?
- 2. An example (James' procedure).
 - a. Logic and theory.
 - b. Examining the data.
 - c. Correlations.
 - d. Finding the best model.
- 3. Evaluating the classical assumptions.

1. What is Regression Analysis?

- It is a statistical method for studying or evaluating the relationship between one or more independent variables and a single dependent variable.
 - We can use regression to predict values of the dependent variable.
 - We do not "prove" causality, we simply test whether the relationships are significantly different from zero.
 - We generally assume the relationships to be linear.
 - Because we cannot explain everything with the independent variables, we include an error term.

Causality



'PLEASE PASS THE SALT'; Examining the Motivational Variables, Idiosyncratic Dynamics And Historic Precedents Associated With the Utterance

BYLINE: By Michael Pacanowsky

SECTION: Outlook; C1

LENGTH: 2561 words

STRONGLY ROOTED in the English speech community is the belief that the utterance, "Please pass the salt," is efficacious in causing salt to move from one end of a table to the source of the utterance. In his "Canterbury Tales," Chaucer notes:

Shee I askked

The salde to passe.

Ne surprissed was I

Tha shee didde (4, p. 318).

Similarly, Dickens writes:

"Old Heep did not become disgruntled at my obstinence. 'Please pass the salt, Davey,' he repeated coldly. I vacillated for a moment longer. Then I passed the salt, just as he knew I would (5, p. 278)."

The question of whether the movement of salt is causally dependent on the utterance of the phrase, "Please pass the salt," has occupied the attention of numerous philosophers (3, 9, 20). Empirical resolution of the validity of this belief, however, was not undertaken until the classic work of Hovland, Lumsdine and Sheffield (8) on the American soldier. Since then, numerous social scientists have explored the antecedent conditions that give rise to this apparent regularity. In this article, we will summarize those efforts that shed some light on the complex phenomenon known as salt passage.

Seminal study upon which all others rest

Army recruits, each sitting alone at one end of a table with salt at the other end, to repeat the utterance, "Please pass the salt," every five minutes for 12 hours. The average distance the salt traveled was .5 inch, which the experimenters explained was due to measurement error. The result of these two studies was, therefore, to demonstrate the importance of the presence of other people in the salt passage phenomenon.

Subsequent research

JANIS AND FESHBACH (10) found that other utterances were just as effective as "Please pass the salt"

No significant differences in the extent of compliance were found due to the utterances, "Please pass the salt," "Would you mind passing the salt?," "Could I have the salt down here, buddy?," and "Salt!" Janis and Feshbach noted that in every successful utterance, the word "salt" was found.

Moderating variables

Festinger (6) tested the effects of substance uncertainty on salt passage. Subjects were placed at a table where salt was loosely piled on a napkin, while sugar was placed in a salt shaker. When a confederate said, "Please pass the salt," the overwhelming number of subjects passed the sugar. From this study, Festinger concluded that the salt shaker, not the salt itself, was the crucial factor in salt passage.

Future research

Third, future research needs to be concerned with the effects of situational variables on salt passage.Kelley's "presence of steak" variable and Milgram's "high threat" variable are suggestive. Effects of information-rich environments, overcrowding, presence of armed conflict, and so on would seem to mediate the salt passage phenomenon.

Simple Regression

 $\mathsf{Y} = \beta_0 + \beta_1 \mathsf{X}_1 + \varepsilon$



Dangers

- Too long model:
 - The model becomes less precise.
- Too short model:
 - The estimated parameters (betas) become biased.
 - The variance in the error term is not correctly estimated.
 - Hypothesis tests and confidence intervals can be misleading.

2a. From Theory ...



 The choice of independent variables should be theoretically based.

 \checkmark theory + logic.

9

... to Equation

 $\mathsf{Y} = \beta_0 + \beta_1 \mathsf{X}_1 + \beta_2 \mathsf{X}_2 + \beta_3 \mathsf{X}_3 + \varepsilon$

- Where:
 - Y = daily earnings,
 - β_0 = the Y intercept,
 - β = the slope coefficients,
 - X = the variables,
 - $\epsilon = error$



Earnings = $\beta_0 + \beta_1$ Seniority + β_2 Gender + β_3 Evaluation + ϵ

2b. Examining the Data

- Sample size.
- The normal distribution and normality.
- Relationships between variables.
- Outliers.
- Missing data.

Sample Size

- Remember that Power = f(effect size, alpha, sample size).
- Effect size = ?
- Alpha RoT = 0.05.
- Sample size = 31
- RoT
- <30 only simple regression.</p>
- 15-20 observations per independent variable.
- With stepwise regression, minimum 50/X variable.

Ramifications of Small Sample

n=31 in our study – so what?

• We are 'overfitting' the variate to the sample.

... Which means...

• We lose generalizability!

Descriptives

IMPORTANT The variables do not need to be normally distributed. The residuals need to be normally distributed.

- Skewness and Kurtosis are OK, but tests of normality indicate not normal.
- IMPORTANT: Box plots indicate that 31 is a consistent outlier.



14

Descriptives continued



Transformation

- Flat Distribution: inverse (1/Y or 1/X).
- Negative Skew: square or cube.
- Positive Skew: log or square root.
- Positive skew means the data is clumped to the left of the histogram.



Transformation

Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
Seniority (in months)	,230	31	,000	,825	31	,000
Senroot	,145	31	,094	,933	31	,052

a. Lilliefors Significance Correction



Positive skew means bunched to the left



Transformed Regression

			Unstandardize	d Coefficients	Standardized Coefficients		
	Model		В	Std. Error	Beta	t	Sig.
	1 (Constant)		70,344	24,746		2,843	,008
	Seniority (in months)		,288	,070	,311	4,089	,000
	Gender		9,769	8,355	,048	1,169	,253
	Firm evaluation of employee		5,858	,651	,698	9,000	,000,
	a. Depende	ent Variable: Daily earning	s				
Interpre	etatior	ר?	Coeff	icients ^a			
			Unstandardize	d Coefficients	Standardized Coefficients		
	Model		В	Std. Error	Beta	l t	Sig.
	1 (C	Constant)	32,638	20,826		1,567	,129
	Se	enroot	5,358	1,626	,258	3,296	,003
	Ge	ender	6,247	8,791	,031	,711	,483
	Fir en	rm evaluation of nployee	6,280	,671	,748	9,364	,000

Coefficients^a

a. Dependent Variable: Daily earnings

2c. Correlations

Low the d	(ROT: significant ependent variab	tly correlated with	All coefficients are the correct sign				
			Daily earnings	Seniority (in months)	Gender	Firm evaluation of employee	
ſ	Daily earnings	Pearson Correlation	1	.900**	.169	.972**	
		Sig. (2-tailed)		.000	.363	.000	
		N	31	31	31	31	
	Seniority (in months)	Pearson Correlation	.900**	1	071	.849**	
		Sig. (2-tailed)	.000		.703	.000	
		Ν	31	31	31	31	
ſ	Gender	Pearson Correlation	.169	071	1	.206	
		Sig. (2-tailed)	.363	.703		.267	
		Ν	31	31	31	31	
ſ	Firm evaluation of	Pearson Correlation	.972**	.849**	.206	1	
	employee	Sig. (2-tailed)	.000	.000	.267		
		Ν	31	/ 31	31	31	

**. Correlation is significant at the 0.01 level (2-tailed).

Danger for multicollinearity (ROT>0.9), therefore, run VIF statistics

Scatterplot



With a correlation coefficient of 0.972 is no surprise that the data are tightly grouped along the line.

Scatterplot – Low Correlation

This is from a different model!

Correlation coefficient = .050



Manufacturer's Image (X4)



Graphs – Chart builder













2d. Finding the Best Model

R² – The coefficient of determination (explained variance)

			Adjusted	Std. Error of
Model	R	R Square	R Square	the Estimate
1	.983 ^a	.966	.962	20.011

Model Summary

a. Predictors: (Constant), Firm evaluation of employee, Gender, Seniority (in months)

Refer to this when discussing individual equations.

Refer to this when comparing equations.

Evaluating the Entire Equation



Use these two when referring to the F-statistic cutoff values.

The F-statistic indicates the significance of the entire equation (ignoring the parts).

F-Distribution $\alpha = 0.05$

Numerator v_1



Evaluating the Parameters (Xs)



Finding the Critical Cutoff Value

- We use n-k-1 to determine the critical value in the table, where
 - N is the number of observations (31),
 - K is the number of parameters (3).
- 31-3-1=27
- The critical value for a one-sided test at a 5% level of significance is 1.703.

t-Distribution



Level of Significance

Rerunning the Model – no gender

Model Summary



a. Predictors: (Constant), Firm evaluation of employee, Seniority (in months)

b. Dependent Variable: Daily earnings

The Final Model



```
Interpreting the Equation
```

Original equation:

 $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \varepsilon$

Final equation:

 $\mathsf{Y} = \beta_0 + \beta_1 \mathsf{X}_1 + \beta_3 \mathsf{X}_3 + \varepsilon$

Which translates to:

```
Earnings = \beta_0 + \beta_1Seniority + \beta_2Evaluation + \varepsilon
```

Or:

```
Earnings = 60.265 + .249Seniority + 6.242Evaluation + ε
Interpretation:
```

Publishing pays!

3. The Classical Assumptions

- This is a fairly basic set of 7 assumptions required to hold for "ordinary least squares" to be considered the best estimator of the regression model.
 - If the assumptions don't hold, other estimation techniques may be better.
 - E.g. Weighted least squares, Two-stage least squares.
 - NOTE: We usually examine standardized residuals to allow for comparison across variables.
 - Studentized are a form of standardized.

The Assumptions

- 1. The regression model is linear in the coefficients.
 - Produce null plot and partial plots.
- 2. Constant variance of the error terms.
 - Plot studentized residuals against predicted values to look for heteroscedasticity.
- 3. Independence of the error terms.
 - Plot residuals against any possible sequencing variable (e.g. Time).
- 4. Normality of the error terms.
 - Test the residuals for a normal distribution (refer to examining data).

The Assumptions

- 5. No explanatory variable is a perfect function of any other explanatory variable.
 - Test for multicollinearity.
- 6. The error term has a zero population mean.
 - This cannot be shown, however it is compensated for by the inclusion of the constant (β₀), so do
 not worry about it.
- 7. The explanatory variables are uncorrelated with the error term.
 - We assume this to be true unless we know otherwise.
 - An example of a violation of this is when introducing an interaction effect as the product of two variables. By definition we have violated this assumption and need to look for an alternative estimation procedure (e.g. 2SLS).





(1) Linearity – Partial Plots



Seniority (in months)

Partial Regression Plot

Dependent Variable: Daily earnings



Firm evaluation of employee



(2) Constant Variance of Error Terms



Hmm?

(3) Independence of the Error Terms

Plot standardized residuals against possible sequencing variables (e.g. time)



(4) Normality of Error Terms

Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
Standardized Residual	,164	30	,038	,915	30	,020

a. Lilliefors Significance Correction







(5) Multicollinearity

- Multicollinearity is a gradual phenomenon so there are no fixed points.
- The correlation coefficient between X_1 and X_2 is 0.849.
 - Rule of thumb: Correlations over .9 indicate a risk of multicollinearity.
- The VIF (variance inflation factor) is 3.582 for each of the variables.
 - Rule of thumb: The VIF should be below 5 (10).

(6) Error has Zero Pop. Mean

- No test
 - Taken care of by the inclusion of the constant (β_0).

(7) X Var. Uncorrelated with Error

Assumed to be true unless otherwise known.

• Example of violation:

•
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon$$

Product
interaction
effect

Advice: look for outliers in partial plots

