



Multiple Regression Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

Analisi Statistica per le Imprese

Prof. L. Neri

Dip. di Economia Politica e Statistica

4.3. Multiple Regression Model



You should be able to:

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

- Understand model building using multiple regression analysis
- Apply multiple regression analysis to business decision-making situations
- Analyze and interpret the computer output for a multiple regression model
- Test the significance of the independent variables in a multiple regression model
- Incorporate qualitative variables into the regression model by using dummy variables



The systematic part

Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

- One may need a mathematical model to quantify the existing relationship between a response variable Y and k explicative variables X_1, \dots, X_k

$$y = f(X_1, \dots, X_k)$$

The multiple linear regression model specifies the functional relationship as

$$f(X_1, \dots, X_k) = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k \quad (1)$$

Geometrically this corresponds to a hyper-plan in k dimensions. The model is extremely useful because:

- ① It has an intuitive geometrical interpretation
- ② It has a simple estimation of the model's parameters



The stochastic part

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

The model always includes a random component, that identifies the stochastic component. This can be expressed as follows:

$$Y = \underbrace{f(X_1, \dots, X_k)}_{\text{systematic}} + \underbrace{\varepsilon}_{\text{stochastic}} \quad (2)$$



The Model specification

Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

Standard notation: for each statistical unit $i=1\dots n$

$$Y_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik} + \varepsilon_i \quad (3)$$

Matrix notation

$$Y = X\beta + \varepsilon \quad (4)$$

Y : $(n \times 1)$ vector of n dependent variable observations

X : $(n \times k)$ matrix of k regressors with each n observations

β : $(k \times 1)$ vector of k parameters

ε : $(n \times 1)$ vector of n



The Matrices

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

$$X = \begin{pmatrix} X_{11} & X_{12} & \cdots & X_{1k} \\ X_{21} & X_{22} & \cdots & X_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{nk} \end{pmatrix}; \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{pmatrix}$$
$$y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}; \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

The matrix X will have a unitary first column if the model is with intercept. In this case the intercept would be β_1 in the multidimensional notation

Standard Assumptions



Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

Assumptions

- The functional relationship must be linear
- Covariates must have a deterministic nature
- The X matrix has full rank
- The error term has a null expected value $E[\varepsilon_i] = 0$
- The error term is homoskedastic: $Var[\varepsilon_i] = \sigma^2$
- The error terms are not correlated: $Cov[\varepsilon_i \varepsilon_j]_{\forall i \neq j} = 0$



OLS estimator

Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

The OLS estimator in multiple linear regression is the vector $\hat{\beta}$ that minimize the following function of $\tilde{\beta}$, where X_i is the i -th row of the X matrix.

$$\min \sum_{i=1}^n e_i^2 = \min \sum (Y_i - X_i \tilde{\beta})^2 \quad (5)$$

$$e = (Y - X \tilde{\beta}) \quad (6)$$

So

$$\tilde{\beta} = (X'X)^{-1} X'Y \equiv \hat{\beta} \quad (7)$$

Notice that the matrix $X'X$ (cross product matrix) must be of full rank in order to be invertible.

OLS estimator properties



Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

$\hat{\beta}$ is BLUE (Best Linear Unbiased Estimator)

One can notice that $\left((X'X)^{-1} X' \right)$ is a matrix of constant elements, therefore $\hat{\beta}$ is a linear transformation of Y .

One can prove that $\hat{\beta}$ is a correct estimator as follows

$$\hat{\beta} = (X'X)^{-1} X'Y = (X'X)^{-1} X'(X\beta + \varepsilon) = \beta + (X'X)^{-1} X'\varepsilon \quad (8)$$

$$E[\hat{\beta}] = \beta + (X'X)^{-1} X'E[\varepsilon] = \beta \quad (9)$$

We could also prove that in the class of the linear and unbiased estimator is the one presenting the minimum variance.



Variance of the OLS estimator

The variance of the OLS estimator is calculated as follows:

$$\text{Var}(\hat{\beta}) = E \left[(\hat{\beta} - \beta) (\hat{\beta} - \beta)' \right] \quad (10)$$

$$\text{Var}(\hat{\beta}) = E \left[(X'X)^{-1} X' \varepsilon \varepsilon' X (X'X)^{-1} \right] \quad (11)$$

$$\text{Var}(\hat{\beta}) = (X'X)^{-1} X' E[\varepsilon \varepsilon'] X (X'X)^{-1} = \sigma^2 (X'X)^{-1} X' X (X'X)^{-1} \quad (12)$$

$$\text{Var}(\hat{\beta}) = \sigma^2 (X'X)^{-1} \quad (13)$$

so for each parameter estimator

$$\text{Var}(\hat{\beta}_j) = \sigma^2 \left[(X'X)^{-1} \right]_{jj} \quad (14)$$



Empirical example

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

A distributor of frozen desert pies wants to evaluate factors that influence demand

- Dependent variable: y = Pie sales (units per week)
- Independent variables: x_1 = Price (\$) and x_2 = Advertising (\$100's)
- Data is collected for 15 weeks

The OLS estimates gives the following estimated model:

$$\hat{y} = 306 - 24x_1 + 74x_2 \quad (15)$$



Interpretation of the estimated coefficient

Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

- each $\hat{\beta}_j$ estimates the average value of Y changes by $\hat{\beta}_j$ units for each 1 unit increase in X_j , holding all other independent variables constant
 - example: $\hat{\beta}_1 = -24$ then sales (y) are expected to decrease, on average, by 24 pies per week for each \$1 increase in selling price (x_1), net of the effects due to advertising (x_2)
 - example: $\hat{\beta}_2 = 74$ then sales (y) is expected to increase, on average, by 74 pies per week for each \$100 increase in advertising (x_2), net of the effects due to price (x_1)
- Intercept
 - the estimated average value of y when all X variables are zero.



Using the model to make predictions

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

We can calculate the predicted sales per week given the selling price (\$5) and advertising expenses (\$350):

$$\hat{y} = 306 - 24(5) + 74(3.5) = 445 \quad (16)$$

We predict to sell 445 pies per week.



The σ^2 estimator

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

$$e = Y - X\hat{\beta} = X\beta + \varepsilon - X \left[(X'X)^{-1} X' (X\beta + \varepsilon) \right] \quad (17)$$

$$e = X\beta + \varepsilon - X\beta - X(X'X)^{-1} X' \varepsilon \quad (18)$$

$$e = \left(1 - X(X'X)^{-1} X' \right) \varepsilon = M_x \varepsilon \quad (19)$$

M_x is a idempotent symmetric matrix. This implies that:

$$M_x = M_x' = M_x^k; \quad \forall k > 0 \quad (20)$$

The σ^2 estimator



Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

$$Q = e'e = \varepsilon' M_x' M_x \varepsilon = \varepsilon' M_x \varepsilon \quad (21)$$

we can prove that

$$E[Q] = \sigma^2 (n - k) \quad (22)$$

so the unbiased estimator of the parameter σ^2 is

$$E\left[\frac{Q}{n-k}\right] = \sigma^2 \Rightarrow \hat{\sigma}^2 = \frac{Q}{n-k} = \frac{e'e}{n-k} \quad (23)$$



ANOVA

Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

Adding to the standard assumptions the following

- The error term has a normal distribution so: $\varepsilon_i \sim N(0, \sigma^2)$

Reminding that

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} = (X'X)^{-1} X'Y \quad (24)$$

follows

$$Y_i = N(X\beta, \sigma^2) \quad (25)$$

$$\hat{\beta}_i = N\left(\beta_i, \sigma^2 \left[(X'X)^{-1} \right]_{ii}\right) \quad (26)$$

$$\hat{\beta} = N\left(\beta, \sigma^2 (X'X)^{-1}\right) \quad (27)$$



ANOVA

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

In order to test the meaning of the whole model, we need to test the hypothesis system

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_1 : \text{almeno un } \beta_j \neq 0$$

the test is

$$F = \frac{ESS/(k-1)}{RSS/(n-k)} = \frac{R^2/(k-1)}{(1-R^2)/(n-k)} \sim F_{(k-1, n-k)}$$

ANOVA

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

	SS	d.f.	Mean Square (MS)
Model	$ESS = \hat{\beta}' X' y = y' y R^2$	k-1	$\hat{\beta}' X' y / (k-1)$
Residual	$RSS = e' e = y' y (1 - R^2)$	n-k	$e' e / (n-k)$
Total	$TSS = y' y = \sum y_i^2$	n-1	

Table 1:

- Construct the F statistic $F = \frac{ESS/(k-1)}{RSS/(n-k)}$
- Find the 95th or the 99th quantile of the distribution $F_{(k-1),(n-k)}$
- If $F > F_{(1-\alpha);(k-1),(n-k)}$ one rejects

 R^2

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

$$R^2 = \frac{ESS}{TSS} \geq 0 \quad (28)$$

$$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{\sum e_i^2}{\sum Y_i^2} \leq 1 \quad (29)$$

$$0 \leq R^2 \leq 1 \quad (30)$$

Adjusted R^2

- R^2 never decreases when a new X variable is added to the model
 - this can be a disadvantage when comparing different models
- What is the net effect of adding a new variable?
 - we lose a degree of freedom when a new X variable is added
 - did the new X variable add enough explanatory power to offset the loss of one degree of freedom?
- The adjusted R^2 adjusts for the number of variables (k).

$$R_{adj}^2 = 1 - \frac{\sum e_i^2 / (n-k)}{\sum Y_i^2 / (n-1)} = 1 - \frac{(n-1)}{(n-k)} (1 + R^2) \quad (31)$$





Empirical example: ANOVA

Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

	SS	d.f.	MS
Model	29460	2	14730
Residual	27033	12	2252
Total	56493	14	

Table 2:

$$R^2 = 29460/56493 = 0.52 \quad (32)$$

52% of the variation in pie sales is explained by the variation in price and advertising

$$R_{adj}^2 = 0.44 \quad (33)$$

44% of the variation in pie sales is explained by the variation in price and advertising, taking into account the sample size and the number of independent variables



Empirical example: is the model significant?

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

- F-test for the overall significance of the model
- Shows if there is a linear relationship between all of the X variables considered together and Y
- use F test statistic
 - in the estimated model the empirical value of F is equal to 6.54
 - the critical value of $F_{0.05;2,12}$ is equal to 3.88
 - $6.54 > 3.88$ therefore the regression model explains a significant portion of the variation in pie sales (there is evidence that at least one independent variable affects y)



Single parameter t-test

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

$$\hat{\beta} \sim N\left(\beta, \sigma^2 \left[(X'X)^{-1} \right]\right) \quad (34)$$

To test the hypothesis if the individual variable X_i has a significant effect on Y we have to test

$H_0: \beta_i = 0$ (no linear relationship)

$H_1: \beta_i \neq 0$ (linear relationship does exist between X_i and Y)

if σ^2 is known, under H_0 :

$$\frac{\hat{\beta}_i - 0}{\sqrt{\sigma^2 \left[(X'X)^{-1} \right]_{ii}}} \sim N(0, 1) \quad (35)$$

Single parameter t-test



Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

Generally σ^2 is unknown and we have to use its estimator

$$\hat{\sigma}^2 = \frac{e'e}{n-k} \quad (36)$$

Therefore the Standard Errors of β is

$$se_{\hat{\beta}} = \sqrt{\hat{\sigma}^2 [(X'X)^{-1}]} \quad (37)$$

Single parameter t-test



Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

under H_0 we have:

$$\frac{\hat{\beta}_i - 0}{\frac{\sqrt{\sigma^2 [(X'X)^{-1}]_{ii}}}{\sqrt{\frac{e'e}{\sigma^2} / (n-k)}}} = \frac{\hat{\beta}_i}{se_{\hat{\beta}}} \sim t_{n-k} \quad (38)$$

- Find the 95th or the 99th quantile of the distribution $t_{(n-k)}$
- if $|t| > t_{(1-\alpha/2);(n-k)}$ one rejects H_0

Empirical example: are individual variables significant?



Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

	coefficients	Standard Error	t
Intercept	306	114.25	2.67
Price	-24	10.83	-2.21
Advertising	74	25.96	2.85

Table 3:

At $\alpha = 0.05$ significant level, the t-value for price is

$|-2.21| > t_{(\alpha/2; 12)} = 2.1788$ so we refuse H_0

At $\alpha = 0.05$ significant level, the t-value for advertising is

$|2.85| > t_{(\alpha/2; 12)} = 2.1788$ so we refuse H_0

There is evidence that both Price and Advertising affect pie sales at $\alpha = 0.05$



Dummy Variables

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

Until now we have assumed that in $Y = X\beta + u$, X are cardinal variables.

One can also use categorial explanatory (i.e. “dummy”) variables that identify specific factors depending on categories:

- Temporal effects
- Spacial effects
- Qualitative variables

We suppose that the Dummies influence just the model intercept (not the slopes).

There will be different regression intercepts corresponding to the different groups/situations, if the dummy variable is significant.



Dummy Variables

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

Lets assume the following generic model:

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 \quad (39)$$

Where Y is the pie sale, X_1 is the price and X_2 is a holiday indicator function, that will be equal to 1 when a holiday has occurred during the week, and equal to 0 when there is no holiday during the week.



Dummy Variables

Multiple Regression Model

Prof. L. Neri

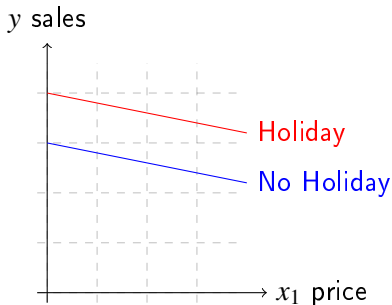
Model, Hypothesis and Estimation

Inference

Dummy Variables

References

$$\begin{aligned}\widehat{Y} &= \beta_0 + \beta_1 X_1 + \beta_2(1) = (\beta_0 + \beta_2) + \beta_1 X_1 \\ \widehat{Y} &= \beta_0 + \beta_1 X_1 + \beta_2(0) = \underbrace{(\beta_0)}_{\text{different intercept}} + \underbrace{\beta_1}_{\text{same slope}} X_1\end{aligned}\quad (40)$$



If $H_0 : \beta_2 = 0$
is rejected
Holiday has
a significant effect
on pie sales



Dummy Variables

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

Example:

$$Sales = 300 - 30(Price) + 15(Holiday) \quad (41)$$

- Sales = number of pies sold per week
- Price = pie price in \$
- Holiday = $\begin{cases} 1 & \text{if holiday has occurred} \\ 0 & \text{if holiday has not occurred} \end{cases}$

As we see the dummy coefficient $\beta_2 = 15$. This implies that on average 15 more pies were sold in weeks with holidays than in weeks without holidays, given the same price.



Dummy Variables

Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References

The number of dummy variables must always be one less than the number of discriminated levels. Imagine that we are analyzing the house market have three different housing levels {ranch, split level, condo}.

Example:

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 \quad (42)$$

Y = house price

X_1 = square meters

$X_2 = \begin{cases} 1 & \text{if ranch} \\ 0 & \text{if not} \end{cases} \Rightarrow \beta_2 \text{ impact of ranch vs. condo}$

$X_3 = \begin{cases} 1 & \text{if split level} \\ 0 & \text{if not} \end{cases} \Rightarrow \beta_3 \text{ split level vs. condo}$



Dummy Variables

Multiple Regression Model

Prof. L. Neri

Model, Hypothesis and Estimation

Inference

Dummy Variables

References

Suppose the estimated equation is:

$$\hat{Y} = 20 + 0.05X_1 + 24X_2 + 15X_3 \quad (43)$$

The we will have as follows:

For a condo ($x_2 = x_3 = 0$):

$$\hat{y} = 20 + 0.05X_1 \quad (44)$$

For a ranch ($x_3 = 0$):

$$\hat{y} = 44 + 0.05X_1 \quad (45)$$

For a split level ($x_2 = 0$):

$$\hat{y} = 35 + 0.05X_1 \quad (46)$$



Multiple
Regression
Model

Prof. L. Neri

Model,
Hypothesis
and
Estimation

Inference

Dummy
Variables

References



Pindyck R. and Rubinfeld D., 1998 “Econometrics Models and Economic Forecasts” (4-th Ed.)



De Luca Amedeo, 1996 “Marketing Bancario e Metodi Statistici Applicati”, vol. 1, Franco Angeli.