



ON SCREEN AND ON AIR TALENT
AN ASSESSMENT OF THE BBC'S APPROACH AND IMPACT
A REPORT FOR THE BBC TRUST

APPENDIX VIII –
THE REGRESSION ANALYSIS

BY OLIVER & OHLBAUM ASSOCIATES

APRIL 2008

APPENDIX VIII – THE REGRESSION ANALYSIS

INTRODUCTION

O&O performed regression analysis to determine the importance of talent on television viewing audiences. The analysis uses two sets of data: a large dataset consisting of strands from the five main networks (BBC1, BBC2, ITV1, Channel 4 and Five) and within sixteen sub-genres¹ and a smaller dataset containing BBC-only data which also included talent costs.

The regression analysis implemented OLS (ordinary least squares) regression on both sets of data, with average audience as the dependent variable (the left-hand side of the equation). The independent variables (right-hand side variables) mainly comprised of dummy variables² to denote strand characteristics, such as the broadcast channel, sub-genre, and stratification of the main talent.

REGRESSION ANALYSIS WITH LARGE DATASET (BBC, ITV, C4 and FIVE)

The equation of the first regression is in log-linear form as denoted below:

$$\ln AUD_i = C + \alpha_1 \text{year06} + \alpha_2 \text{year07} + \beta_1 \text{BBC2} + \beta_2 \text{ITV} + \beta_3 \text{CH4} + \beta_4 \text{FIVE} + \beta_5 \text{chat} + \beta_6 \text{family} + \beta_7 \text{othercomedy} + \beta_8 \text{quiz} + \beta_9 \text{sitcom} + \beta_{10} \text{science} + \beta_{11} \text{naturalhist} + \beta_{12} \text{serials} + \beta_{13} \text{series} + \beta_{14} \text{plays} + \beta_{15} \text{factent} + \beta_{16} \text{human} + \beta_{17} \text{factdrama} + \beta_{18} \text{lifestyle} + \beta_{19} \text{cook} + \beta_{20} \text{prepeak} + \beta_{21} \text{earlypeak} + \beta_{22} \text{midpeak} + \beta_{23} \text{latepeak} + \beta_{24} \text{afterpeak} + \beta_{25} \text{SUPER} + \beta_{26} \text{TOP} + \beta_{27} \text{KNOWN} + \varepsilon$$

Where the variables are defined as the following:

VARIABLE	DEFINITION
<i>C</i>	Constant
<i>lnAUD</i>	Natural log of viewing audience
<i>year06</i>	Dummy variable takes the value of 1 if strand occurred in 2006
<i>year07</i>	Dummy variable takes the value of 1 if strand occurred in 2007
<i>BBC2</i>	Dummy variable takes the value of 1 if strand occurred on BBC2
<i>ITV</i>	Dummy variable takes the value of 1 if strand occurred on ITV1
<i>CH4</i>	Dummy variable takes the value of 1 if strand occurred on Channel 4
<i>FIVE</i>	Dummy variable takes the value of 1 if strand occurred on Channel 5
<i>chat</i>	Dummy variable takes the value of 1 if strand fell into the chat sub-genre
<i>family</i>	Dummy variable takes the value of 1 if strand fell into the family show sub-genre
<i>othercomedy</i>	Dummy variable takes the value of 1 if strand fell into the other comedy sub-genre
<i>quiz</i>	Dummy variable takes the value of 1 if strand fell into the quiz/panel game sub-genre
<i>sitcom</i>	Dummy variable takes the value of 1 if strand fell into the sitcom sub-genre

¹ The sixteen sub-genres are: Chat shows, family shows, quiz/panel games, other comedy, sitcoms, drama series, drama serials, drama single plays, factual entertainment, human interest, factual drama, lifestyle, and cookery.

² Dummy variables are binary and take the value of "1" when the status is "true" or "0" when the status is "not true." For example, the dummy variable *ITV* would take the value of "1" if the strand was broadcast on ITV1.

<i>science</i>	Dummy variable takes the value of 1 if strand fell into the science/medical sub-genre
<i>naturalhist</i>	Dummy variable takes the value of 1 if strand fell into the natural history sub-genre
<i>serials</i>	Dummy variable takes the value of 1 if strand fell into the drama-serials sub-genre
<i>series</i>	Dummy variable takes the value of 1 if strand fell into the drama-series sub-genre
<i>plays</i>	Dummy variable takes the value of 1 if strand fell into the drama-single plays sub-genre
<i>factent</i>	Dummy variable takes the value of 1 if strand fell into the factual entertainment sub-genre
<i>human</i>	Dummy variable takes the value of 1 if strand fell into the human interest sub-genre
<i>factdrama</i>	Dummy variable takes the value of 1 if strand fell into the factual drama sub-genre
<i>lifestyle</i>	Dummy variable takes the value of 1 if strand fell into the lifestyles and makeover sub-genre
<i>cook</i>	Dummy variable takes the value of 1 if strand fell into the cookery show sub-genre
<i>prepeak</i>	Dummy variable takes the value of 1 if strand occurred between 4pm-6pm
<i>earlypeak</i>	Dummy variable takes the value of 1 if strand occurred between 6pm-8pm
<i>midpeak</i>	Dummy variable takes the value of 1 if strand occurred between 8pm-10pm
<i>latepeak</i>	Dummy variable takes the value of 1 if strand occurred between 10pm-11pm
<i>afterpeak</i>	Dummy variable takes the value of 1 if strand occurred between 11pm-6am
<i>SUPER</i>	Dummy variable takes the value of 1 if strand had "super talent"
<i>TOP</i>	Dummy variable takes the value of 1 if strand had "top talent"
<i>KNOWN</i>	Dummy variable takes the value of 1 if strand had "known talent"

The variable coefficients (the β 's) measure the influence that variable has on the dependent variable. However, when a regression equation is in log-linear form, a calculation is required to derive the true influence of the dummy variables.

If the general equation takes the form $\ln Y = \beta_1 + \beta_2 X + \beta_3 D_1$ where D is a dummy variable, the percentage increase in Y for the dummy taking a value of "1" compared to the base case is $(e^{\beta_3} - 1) * 100$

REGRESSION RESULTS

The results from the first regression are presented below.

Dependent Variable: LN_AUD	COEFFICIENTS	STD ERROR	T-STAT	SIGNIFICANCE	DUMMY INFLUENCE (%)
<i>Constant</i>	5.45	0.16	33.48	0.00	
<i>BBC2</i>	-0.63	0.33	-1.91	-0.06	<i>-46.7</i>
<i>ITV1</i>	-0.18	0.07	-2.41	0.02	<i>-16.2</i>
<i>CH4</i>	0.18	0.17	1.06	0.29	<i>19.4</i>
<i>FIVE</i>	-0.28	0.12	-2.27	0.02	<i>-24.2</i>
<i>YEAR 2006</i>	-0.16	0.06	-2.63	0.01	<i>-14.6</i>
<i>YEAR 2007</i>	-0.27	0.06	-4.22	0.00	<i>-23.9</i>
<i>CHAT SHOWS</i>	0.18	0.24	0.73	0.46	<i>19.2</i>
<i>FAMILY SHOWS</i>	1.16	0.19	6.27	0.00	<i>220.6</i>
<i>OTHER COMEDY</i>	0.31	0.21	1.48	0.14	<i>36.7</i>
<i>QUIZ/PANEL SHOWS</i>	0.64	0.20	3.19	0.00	<i>89.7</i>
<i>SITCOMS</i>	0.23	0.21	1.09	0.28	<i>25.6</i>
<i>SCIENCE/MEDICAL</i>	0.40	0.25	1.59	0.11	<i>48.6</i>
<i>NATURAL HISTORY</i>	0.50	0.20	2.50	0.01	<i>65.5</i>
<i>SERIALS</i>	1.06	0.22	4.79	0.00	<i>188.8</i>
<i>SERIES</i>	1.12	0.19	5.96	0.00	<i>206.9</i>
<i>SINGLE PLAYS</i>	1.06	0.20	5.43	0.00	<i>188.6</i>
<i>FACTUAL ENTERTAINMENT</i>	1.12	0.17	6.55	0.00	<i>207.9</i>
<i>HUMAN INTEREST</i>	0.42	0.15	2.79	0.01	<i>52.1</i>
<i>FACTUAL DRAMA</i>	1.36	0.42	3.26	0.00	<i>289.2</i>
<i>LIFESTYLE</i>	0.03	0.16	0.15	0.88	<i>2.6</i>
<i>COOKERY</i>	-0.24	0.21	-1.18	0.24	<i>-21.7</i>
<i>PRE PEAK</i>	0.26	0.15	1.68	0.09	<i>29.2</i>
<i>EARLY PEAK</i>	0.56	0.10	5.71	0.00	<i>75.2</i>
<i>MID PEAK</i>	0.92	0.10	9.45	0.00	<i>151.4</i>
<i>LATE PEAK</i>	0.17	0.12	1.37	0.17	<i>18.2</i>
<i>AFTER PEAK</i>	-0.60	0.13	-4.67	0.00	<i>-45.2</i>
<i>SUPER TALENT</i>	1.33	0.10	12.89	0.00	<i>278.9</i>
<i>TOP TALENT</i>	1.23	0.09	13.16	0.00	<i>243.5</i>
<i>KNOWN TALENT</i>	0.98	0.09	10.98	0.00	<i>167.3</i>
R²=66%					

Variables in bold indicate that they are 95% significant which means that there is only a 5% chance that their influence would occur randomly.

Firstly, since the equation is in log linear form, the correct interpretation of the dummy variable coefficients is derived from the calculation stated in the equation above and denoted in the last column of the table. Secondly, dummy variables can only be interpreted relative to the base case. In regression equations the base case is implicit. For example in the above equation, the base sub-genre is 'history.' This means that all

sub-genres have a representative dummy variable in the equation *except* 'history' since 'history' would be the implied sub-genre if all sub-genre dummy variables took the value of "0." In the regression above, the base time slot is 'daytime', the base channel is 'BBC1', the base year is '2005' and the base level of talent is 'unknown.' So for example, with all else being equal, a strand with 'super talent' on average has 278% higher audiences relative to shows without any talent.

MODEL LIMITATIONS:

- The large dataset is still a very small dataset when compared to the entire population of programmes on the five networks. Not all genres are represented and of the genres that are in the model, only a subset of strands within those subgenres is included. Including only a small proportion of the entire population of data possibly introduces a selection bias, although it is unknown if this has actually occurred in this case.
- One-off programmes and repeats are not removed from this dataset, and there is a high likelihood of substantial programme repeats contained within the dataset. Additionally, the BARB average audience measure will also contain repeats, in other words the average audience of a strand will be influenced by audiences viewing its repeat showing. Programmes with a higher repeat rate will be more susceptible to this type of data bias. This bias may especially present a problem where shows are repeated in the daytime or during late nights, but where the original show was broadcast during peak time. Thus one could attribute some degree of measurement error to the audience viewing data.
- Where appropriate, the talent value survey provided influence on the coding of talent. However, the majority of talent coding had to be done on a subjective basis. There is a potential the data to be miscoded for some strands.
- There has been no testing done with respect to the run of causality, i.e. do large audiences watch programmes and therefore create talent (bigger audiences lead to higher talent costs because channels can afford it) or does talent drive larger audiences. This type of testing (Granger causality testing) has not been performed in this analysis.
- A potentially more robust regression analysis would involve not only the entire population of strands but also data over time. A time series approach is more insightful because it facilitates the use of differenced equations which would capture the *changes* in talent over time. Secondly, a differenced equation would eliminate unobserved effects not captured by the measured variables. There exists the possibility that some effects are unobservable either because there exists no tangible data in which to measure the effect or simply because there is no access to the data which would capture it. If it is the case that these 'un-captured' effects are correlated or enveloped within another variable that is included in the regression, then the coefficient of the variable will not distinguish between two potentially distinct effects. Thus the significance and measure of the coefficient for the variable will be incorrect. Lastly, a time series approach would also enable the representation of 'strand time-specific' effects, e.g. legacy, loyalty, cult status, etc.

There exists the potential for multi-collinearity: when the variables on the right-hand side of the equation are correlated. For example, some genres might only be shown in peak times so the timeslot and genre might correlate with each other. This will also have an impact of the variable coefficients and may make them seem significant when in actuality they are not.

REGRESSION ANALYSIS WITH BBC DATASET (SMALL BBC SET)

A second regression analysis used a much smaller data set containing limited BBC-only information. The regression equation is also log-linear and takes the following form:

$$\ln AUD_i = C + \alpha_1 year06 + \alpha_2 year07 + \beta_1 BBC2 + \beta_2 BBC3 + \beta_3 earlypeak + \beta_4 midpeak + \beta_5 serials + \beta_6 plays + \beta_7 panel + \beta_8 family + \beta_9 sitcom + \beta_{10} othercom + \beta_{11} lifestyle + \beta_{12} natural + \beta_{13} history + \beta_{14} human + \beta_{15} \ln talent\ cost + \varepsilon$$

Where the variables are defined as the following:

VARIABLE	DEFINITION
<i>lnAUD</i>	Natural log of viewing audience
<i>year06</i>	Dummy variable takes the value of 1 if strand occurred in 2006
<i>year07</i>	Dummy variable takes the value of 1 if strand occurred in 2007
<i>BBC2</i>	Dummy variable takes the value of 1 if strand occurred on BBC2
<i>BBC3</i>	Dummy variable takes the value of 1 if strand occurred on BBC3
<i>earlypeak</i>	Dummy variable takes the value of 1 if strand occurred between 6pm-8pm
<i>midpeak</i>	Dummy variable takes the value of 1 if strand occurred between 8pm-10pm
<i>serials</i>	Dummy variable takes the value of 1 if strand belongs to the drama: serials sub-genre
<i>plays</i>	Dummy variable takes the value of 1 if strand belongs to the drama: plays sub-genre
<i>panel</i>	Dummy variable takes the value of 1 if strand belongs to the quiz/panel game sub-genre
<i>family</i>	Dummy variable takes the value of 1 if strand belongs to the family show sub-genre
<i>sitcom</i>	Dummy variable takes the value of 1 if strand belongs to the sitcom sub-genre
<i>othercom</i>	Dummy variable takes the value of 1 if strand belongs to the other comedy sub-genre
<i>lifestyle</i>	Dummy variable takes the value of 1 if strand belongs to the lifestyle sub-genre
<i>natural</i>	Dummy variable takes the value of 1 if strand belongs to the natural history sub-genre
<i>history</i>	Dummy variable takes the value of 1 if strand belongs to the history sub-genre
<i>human</i>	Dummy variable takes the value of 1 if strand belongs to the human interest sub-genre
<i>ln talentcosts</i>	The natural log of talent costs for the strand that year

This regression equation is also log linear but in this case the talent measure used is not a dummy variable, but rather the actual talent costs for the strand. Thus one can interpret the coefficient on the talent variable without manipulation. So if the talent cost coefficient is 0.2, then a 1% increase in talent costs increases the audience by 0.2%. The

other variables are dummy variables and their coefficients should be interpreted as in the first regression.

Dependent Variable: LN AUD	COEFFICIENT	STD. ERROR	T-STAT	SIGNIFICANCE	DUMMY SIGNIFICANCE (%)
(Constant)	4.84	1.27	3.81	0.00	
YEAR 06	-0.13	0.12	-1.09	0.28	-11.8
YEAR 07	-0.14	0.12	-1.17	0.25	-12.8
EARLY_PEAK	0.68	0.29	2.35	0.02	96.5
MID_PEAK	0.48	0.17	2.83	0.01	62.2
BBC2	-0.68	0.16	-4.16	0.00	-49.4
BBC3	-1.88	0.39	-4.85	0.00	-84.7
SERIALS	0.07	0.21	0.32	0.75	6.8
PLAYS	0.23	0.34	0.66	0.51	25.3
PANEL GAMES	-0.23	0.47	-0.48	0.63	-20.2
FAMILY SHOWS	0.49	0.16	3.04	0.00	63.6
SITCOM	-0.38	0.27	-1.41	0.16	-31.8
OTHER COMEDY	-0.01	0.31	-0.02	0.98	-0.7
NATURAL HISTORY	0.94	0.38	2.49	0.02	157.0
LIFESTYLE	0.22	0.26	0.85	0.40	24.4
HISTORY	0.66	0.40	1.65	0.10	94.4
HUMAN INTEREST	0.50	0.47	1.05	0.30	64.5
LN_COST	0.21	0.09	2.45	0.02	23.9
R² =91%					

In this regression the base variables were 'BBC1' for the channel, 'off-peak' for the time-slot, 'drama series' for sub-genre, and '2005' for year. Again, variables that are 95% significant are denoted in bold type.

MODEL LIMITATIONS:

- The BBC dataset was a very small dataset (less than 100 observations and less than 50 strands) which limits meaningful extrapolations. Because the sample is small, there may not be sufficient variety represented by the dataset and there is a very strong possibility of selection bias in this sample. For example, one cookery strand would not be indicative of the entire cookery sub-genre, but using only one strand within a regression analysis would give that strand undue influence.
- The talent cost data represents more than just the cost of actors and presenters. This would mask or over-emphasise the true influence of pure talent. If this measurement error is combined with the limitations of a small dataset, the bias is magnified.