

Statistical Analysis Report

s.to

26 November 2020



000099





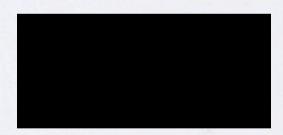
Contents

Statement of	3
Executive Summary	4
Data Collection	4
Statistical Sampling and Analysis	4
Results	6



Statement of

- I have been instructed by the Motion Picture Association to provide a statistical analysis of the proportion
 of copyright protected content on the website identified in the summary table in the 'Results' section of
 this report (the Target Website).
- 2. I work as a Data Scientist and Statistician for INCOPRO Limited (INCOPRO) and have held this position since 18 December 2014. INCOPRO is a limited company incorporated in England and Wales which monitors the online environment for infringement of intellectual property rights. INCOPRO specialises in the identification of infringement of intellectual property and in the provision of reporting and other services to assist in the evaluation of action required to protect intellectual property online. Part of my role within INCOPRO is to perform statistical analysis and to develop systems to apply a statistical approach to the collection of data.
- 3. I have 17 years' professional experience in various roles with a strong quantitative and computational element, including 5 years as an interdisciplinary Post-Doctoral Research Associate at Imperial College London specialising in the mathematical modelling and statistical analysis of large data sets. I hold a BSc in Biological Sciences from Birkbeck College London, an MRes in Modelling Biological Complexity and Bioinformatics from University College London (UCL) and a Ph.D. in Mathematical Modelling from the University of London. I am a Fellow of the Royal Statistical Society.
- In order to enable the analysis summarised in this report, I was supplied with data in the form of a website sample scrape undertaken by Incopro's Analysts (the Analysts).
- 5. I confirm that I have made clear which facts and matters referred to in this report are within my own knowledge and which are not. Those that are within my own knowledge I confirm to be true. The opinions I have expressed represent my true and complete professional opinions on the matters to which they refer.



Date: 26 November 2020



Executive Summary

- 6. I analysed the on the Target Website using recognised methods of statistical sampling.
- My conclusion is that

Data Collection

- On 10 November 2020, the Analysts used random sampling to scrape (or copy) 100 links to available TV content on the Target Website.
- The Analysts produced the scraped sample, set out as a computer text file, with each line in the file representing a link to relevant content. Hereafter, I will refer to the computer file as the "sample".

Statistical Sampling and Analysis

- Inferential statistics permits the drawing of conclusions as to the characteristics of the whole population based on the results of a random sample drawn from the whole population. By undertaking this sampling procedure, I wanted to be able to draw conclusions as to made available by the Target Website based on a sample drawn from each of them. Taking a random sample ensures a good representative sample because it removes any human influence, and therefore subjectivity or bias from the process. In contrast, any deterministic method, however arbitrary or devoid of subjectivity, would run the risk of returning biased results. To evaluate on the Target Website a random sample of the entries was drawn from each of them to be analysed manually.
- The simplest method of drawing a sample from a population is sampling with replacement. This method was implemented in respect of all Target Website.
- Sampling with replacement means that individual records from the population might appear more than once in the sample. This method ensures that the precision of the results depends on the size of the sample only. In contrast and counter-intuitively, the sampling rate (defined as the ratio of the sample size to the population size) plays no role in that precision. This is because sampling with replacement ensures that each time a record is drawn from the population, the probability that remains constant throughout the process. In other words, the sampling rate is immaterial because the underlying population size is immaterial.
- In essence, the process replicates a series of random experiments with a binary outcome and recording the overall result. The canonical example of such a random experiment is throwing a coin and checking on which side it lands. The proportion of heads obtained by throwing a well-balanced coin is more likely to be



close to ½ with a large number of throws. Likewise, one is less likely to obtain an unbroken series of heads (or tails, for that matter) if that coin is thrown a large number of times. By sampling randomly with replacement one aims to replicate this kind of property and thereby harness the laws of probability. Of course, the "coin" that is virtually thrown might be unbalanced and the degree to which it is unbalanced is investigated through sampling.

The sample size of 100 was selected on the basis that this provides enough precision to the results to be able to confidently reach robust conclusions as to the proportion of content that is commercially available or otherwise protected by copyright made available by the Target Website.
 Subsequently, the Analysts manually analysed the 100 selected entries for the Target Website to determine the number of entries that related

The Analysts used the following methodology:

· They analysed every record starting with the first.

	when the sales of a street	d also	A COURT TO SERVICE	
or each e	entry, they checke	d the		
Where the	ey were unable to	identify an ent	ry	

- The list of content titles contained within the sample and analysed using the above method can be found in Appendix 1.0: List of Titles.
- 17. The results of this analysis enabled me to make inferences about on the Target Website as a whole.
- 18. To reflect the fact that the entire population has not been analysed, the results of inferential statistics are typically reported by reference to a degree of confidence and/or precision: it is not valid to just return a single indicator,
 It needs to be accompanied with complementary indicators describing the confidence and/or

precision. Here, I have chosen to present the results in two different ways.



- 19. Firstly, by reference to a 99% confidence interval, which is the traditional way to report results from surveys and polls. This means that instead of a single value, a range of values is returned and there is a 99% chance that the corresponding population value is within this range.²
- 20. Secondly, by contrasting the probability that the population proportion is above a certain value to the probability that it is below that value. This allows me to assert, more directly, where the preponderance of evidence lies. In order to do this I define the following descriptive proportions:
 - A 'large majority': 90%.
 - A 'sizeable majority': 80%.
 - A 'simple majority': 50%.
- 21. The obvious consequence of those definitions is that asserting that the large majority of the examined content of a website is infringing is a stronger statement than asserting that a sizeable majority is, which in turn is a stronger statement than asserting that a simple majority is.
- 22. It is important to note that the results obtained are not absolute because only a portion of the underlying population was analysed. The robust statistical methodology employed here is designed to mitigate as far as possible the likelihood of an erroneous result.
- 23. Both of the series of indicators were computed using another bespoke computer program written in "R", a computer language and environment widely used and recognised to perform statistical analysis, which has become a standard tool within the statistical community and with which I am very familiar. The raw results that the Analysts communicated to me were fed to that program to obtain the indicators presented below.
- 24. Within that program, I used specialised commands existing within the R environment implementing the calculations of the aforementioned indicators.³ The mathematical formulas involved in these calculations are quite complicated and would otherwise have taken a long time to deploy using less specialised tools. I performed so-called sanity checks to verify the soundness of the computer program's output by using approximate formulas for the confidence intervals and comparing these results to the ones returned by the program. These checks were performed before the actual data presented below was fed to the program and were successful, which convinced me that the program could be confidently used with real data.

Results

25. In respect of s.to I was informed by the Analysts that they were able to match 100 entries from the analysed sample of 100 entries as

On that basis, I conclude with 99% confidence

comprised between

In addition, the probability that this very close to 100% whereas the

² Strictly speaking that probability of 99% refers to the range itself rather than the population value: there is a 99% chance that the range comprises the real value.

³ The confidence intervals were computed using the "exact" method in the R binom package. The second set of indicators was determined using Bayesian inference with a uniform prior on the interval [0,1]. The prior and posterior distribution was modelled with the Beta distribution as outlined in *Doing Bayesian Analysis* by John K. Kruschke (Chapter 6).

000104





	probability that	is very small:	54. TH	herefore the prepare	onderance of evidence
	overwhelmingly lies towards the f	act that		on s.to	
26.	The table below collates the resu	Its of the sample ana	alvsis for th	ne Target Website.	The preponderance of
	evidence lies towards the fact th		***************************************		f the Target Website is
				overwhelmi	ngly so in this case.

⁴ In order to represent a very small number it is convenient and customary to use either scientific or E notation where either an 'E' or an 'e' represents "times ten raised to the power of". In this case, which is words would be times 10 raised to the power of -3". Converting to standard notation would therefore equal 0.00239% (there are 2 zeros following the decimal point and preceding the standard notation).

Farget Website	Categories Scraped	Sample size	Sample No. of links size matched to content that is commercially available	Lower bound of 99% confidence interval	Upper bound of 99% confidence interval	Definition of "the majority"	ority"
	VI	100	100	94.84%	100%	%06	

1																											1
	No.	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	20	
	No.	1	2	3	4	5	9	7	89	6	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	

Appendix 1.0: List of Titles

	No. of the last of										THE RESERVE OF THE							Mesen The control of							
No.	76	77	78	79	80	81	82	83	84	85	86	87	88	89	06	91	92	93	94	95	96	97	98	66	100
										100 mm															
No.	51	52	53	54	55	95	57	58	59	09	61	62	63	64	65	99	29	89	69	70	71	72	73	74	75



