



No. 11-2021

Maximilian Maurice Gail and Phil-Adrian Klotz

**The Impact of the Agency Model on E-book Prices:
Evidence from the UK**

This paper can be downloaded from:

<https://www.uni-marburg.de/en/fb02/research-groups/economics/macroeconomics/research/magks-joint-discussion-papers-in-economics>

Coordination: Bernd Hayo • Philipps-University Marburg
School of Business and Economics • Universitätsstraße 24, D-35032 Marburg
Tel: +49-6421-2823091, Fax: +49-6421-2823088, e-mail: hayo@wiwi.uni-marburg.de

The Impact of the Agency Model on E-book Prices: Evidence from the UK

Maximilian Maurice Gail^a, Phil-Adrian Klotz^{a,*}

^a*Chair for Industrial Organization, Regulation and Antitrust, Department of
Economics, Licher Strasse 62, Giessen, 35394, Germany*

Abstract

This paper empirically analyzes the effect of the widely used agency model on retail prices of e-books sold in the United Kingdom. Using an unique cross-sectional dataset of e-book prices for a large number of book titles across all major publishing houses, we exploit cross-genre and cross-publisher variation to identify the effect of the agency model on e-book prices. Since the genre information is ambiguous and even missing for some titles in our original dataset, we also apply a Latent Dirichlet Allocation (LDA) approach to determine detailed book genres based on the book's descriptions. We find that retail prices for e-books sold under the agency model are on average 18% cheaper than book titles sold under the wholesale model. Our results are robust to different regression specifications, an instrumental variable approach, and double machine learning techniques.

Keywords: e-books, agency agreements, vertical restraints, Amazon, double machine learning

JEL: D12, D22, L42, L81, L82, Z11

Declarations of interest: None

1. Introduction

The rise of the internet - and platform markets more specifically - has accelerated the usage of so-called agency arrangements. Thereby, suppliers

*Corresponding author

Email address: phil.a.klotz@wirtschaft.uni-giessen.de (Phil-Adrian Klotz)

Preprint submitted to Elsevier

May 23, 2022

pay retailers sales royalties to distribute products at prices determined by suppliers. On the contrary, in many traditional retail markets suppliers charge retailers wholesale prices and retailers set final prices to consumers (the "wholesale" model). Even though agency arrangements are also used in some of these conventional markets (e.g., newspapers sold at kiosks), this form of a vertical contract is especially prevalent in online markets (e.g., Amazon Marketplace, Apple App Store, eBay Buy It Now). It is also frequently used in the retail market examined in this paper, namely the e-book market.

Books are experience goods because readers can ascertain the quality only after reading a given book title (Nelson, 1970; Reimers and Waldfogel, 2021). In some countries such as Germany, France or Japan, book prices are fixed whereas countries without fixed book price (FBP) systems include the UK and the USA. Fixed book prices are a form of resale price maintenance (RPM) where publishers set retail prices and price competition between retailers is restricted or completely eliminated.¹ Mostly the motivation behind the introduction of FBP systems is the assurance for a broad and diverse supply of books, available through a geographically wide network of bookstores.

With the advent of e-books, countries with a FBP system for print books had to decide whether to extend existing legislation to e-books. It is ques-

¹Presently, 15 OECD countries have a regulation for fixing the prices of printed books. The fixed prices for printed books typically last 18-24 months after a book has been published.

tionable whether the same cultural policy arguments and legal considerations apply, in particular because a geographically wide network of bookstores is irrelevant for e-books. Nevertheless, eight OECD countries with fixed prices for printed books also have fixed prices for e-books, while no country is known to have such a RPM regulation for e-books but not for print books (Poort and van Eijk, 2017). However, in many countries without fixed prices for e-books (such as the UK) this digital product is partly sold under the agency model (Gilbert, 2015), which has similar effects as RPM between a manufacturer and a retailer.²

With regard to e-books, Apple in co-operation with the six largest publishing houses have been the first who adopted the agency model in 2010 in response to Amazon's aggressive pricing strategy to gain market share. In April 2012, the Department of Justice (DOJ) sued Apple and five of the six publishing houses for conspiring to raise e-book prices by using the agency model in conjunction with most-favored nation (MFN) clauses.³ Three of the publishers settled shortly after the antitrust case was filed, while the other two followed later the same year, which meant that the five publishers

²They are economically similar in the sense that the upstream firms control the retail price. However, a major difference is that for the agency pricing the downstream firms individually delegate retail pricing to the upstream firms, whereas under the classical case RPM is imposed at the market level for any given good.

³ See *United States v. Apple Inc.*, 12 Civ. 2826 (DLC). MFN clauses stipulate that the retail price set by a publisher through one retailer can be no higher than the retail price set by that publisher through a competing retailer. Hence, MFN guarantees a retailer who prefers a higher commission, if it raises the commission it charges for one publisher, the retail price will remain the same relative to the other retailers. This effect encourages retailers to push for higher fees, which results in higher retail prices (Johnson, 2017).

could not restrict a retailer's ability to set e-book prices for a period of two years.

Empirical evidence on the price effects of RPM and fixed book prices as well as of the agency model is scarce. While systematic empirical evidence on RPM is limited to case studies ([MacKay and Smith, 2017](#); [Ippolito, 1991](#)), the only study investigating the empirical effect of the agency model is the one from ([De los Santos and Wildenbeest, 2017](#)). They have used data on e-book prices of bestselling book titles for the years 2012 and 2013 and the Apple case as an exogenous shock to show that the agency model in combination with MFN clauses led to an average increase in prices between 8-18% (depending on the retailer).

The goal of this study is to analyze the price effect of the agency model using a larger and more detailed data set (especially not only incorporating bestselling book titles) and to check whether similar effects also occur in the absence of an alleged conspiracy as in the Apple case. The internet allows consumers access to a larger number of book titles rather than simply the popular ones (bestselling titles) and different authors have shown the importance of those "long tail" titles in markets for creative products (e.g., [Aguiar and Waldfogel, 2018](#); [Brynjolfsson et al., 2003](#)).⁴

Our cross-sectional dataset contains prices for 12,001 e-books published on Amazon UK between 2010 and 2020. Using data from Amazon ensures a high market coverage since Amazon accounted for 50 percent of the UK

⁴[Brynjolfsson et al. \(2003\)](#) have estimated that the benefit consumers obtain from access to long tail book titles may be as high as \$1.03 billion per year alone in 2000.

book sales in 2018.⁵ We further use publisher- and book genre variation to estimate the effect of the agency model on e-book retail prices.

The results of our propensity score matching design indicate that e-books sold under the agency model on *Amazon.co.uk* are on average 18% cheaper than digital books sold under the wholesale model. An instrumental variable approach and various double machine learning (DML) concepts support our main finding. This result contradicts the empirical outcome from [De los Santos and Wildenbeest \(2017\)](#)⁶, but fits into explanations put forward by the theoretical literature on agency versus wholesale models ([Johnson, 2020](#); [Foros et al., 2017](#); [Gaudin and White, 2014](#)).

The rest of the paper is structured as follows. In Section 2, we describe the related literature. We present our unique data set in Section 3.1. Descriptive statistics are given in Section 3.2 and our text mining approach to determine book genres is explained in Section 3.3. Section 4 presents our main estimation strategy and results. In Section 5, our robustness checks are outlined. We conclude and outline the contributions of our paper in

⁵ See Nielsen (2018), "Books & Consumers - UK Industry Standard Report Q4 2018", p. 13.

⁶Our paper differs from the study of [De los Santos and Wildenbeest \(2017\)](#) in three important ways. First, while [De los Santos and Wildenbeest \(2017\)](#) use the court decision in the Apple Case (see Footnote 3) as an exogenous shock in their approach, our findings do not rely on an alleged conspiracy. Second, we not only incorporate bestselling book titles into our empirical analysis but also books from the "long tail". And third, [De los Santos and Wildenbeest \(2017\)](#) cannot measure the "pure" price effect stemming from agency agreements since those were used in conjunction with MFN clauses at that time. Our price effect can be attributed only to the agency agreements because Amazon has settled with the EU Commission in 2017 not to include MFN clauses in respect of any e-book distributed in the EEA for the next five years (See AT.40153 E-book MFNs and related matters (Amazon), Decision dated May 4, 2017. We have scraped the Amazon data in 2020).

Section 6.

2. Related Literature

Our article contributes to several strands of literature. First and foremost, it is related to studies which investigate the competitive effects of the agency model. While the empirical literature on the economic effects of the agency model is rather scarce (an exception is the study of [De los Santos and Wildenbeest, 2017](#)), several recent theoretical papers have analyzed differences in retail prices between the agency and the wholesale model. One strand of this literature is focused on a lock-in effect of consumers in this context. [Johnson \(2020\)](#) finds that when publishers set retail prices instead of retailers (agency model), prices may be higher in early periods but lower in later periods since in the wholesale model retailers initially set low prices to lock in consumers, but find it optimal to raise prices once a sufficient number of consumers are locked in.

Another strand of the theoretical literature on agency models assumes that complementary devices are necessary for the enjoyment of the main products (e.g., an e-book reader in the case of e-books). [Gaudin and White \(2014\)](#) point out that the incentive of a retailer to set high prices is higher when she has monopolistic control over a complementary device, as it was the case in the e-book market when e-books from Amazon could only be read on a Kindle device. In another model-theoretical setup, [Abhishek et al. \(2016\)](#) show that agency selling is more efficient than the wholesale model and leads to lower retail prices, even though retail prices may be

higher under the agency model if there are positive externalities from sales of associated products (such as e-readers in the case of e-books).

[Foros et al. \(2017\)](#) show that the agency model is always anti-competitive (leads to higher retail prices) when it is adopted by the platforms on a market-by-market basis. To be more specific, they find that upstream firms (publishers) will set higher retail prices than downstream firms (retailers) would set if they were in control as long as competition is greater among retailers than among publishers. Moreover, they point out that a retailer who sets retail prices independently (wholesale model) benefits when a horizontal rival is restricted by the agency model since the latter creates a price umbrella which makes it profitable for the independent price-setting retailer to increase prices. [Condorelli et al. \(2018\)](#) present a theory that makes the decision whether to use agency or wholesale models endogenously in an environment where the retailer has privileged information about the valuations of consumers and show that retailers prefer the agency model.

More generally, our article is also related to the broader literature on RPM. RPM can lead to lower retail prices because of the internalization of vertical externalities such as double marginalization. Since a manufacturer chooses the wholesale price given its costs and the retailer as a result chooses the retail price given the wholesale price, both firms will add their mark-up and in the end the consumers will pay too high prices. An obvious possibility to solve this problem is RPM because the manufacturer simply imposes the resale price on the retailer ([Spengler, 1950](#); [Tirole, 1988](#)).

Moreover, [Telser \(1960\)](#) and [Yamey \(1954\)](#) were the first who noted that strong intra-brand competition can be detrimental to retailer's incentives to invest in free-rideable services. Suppose a situation in which a retailer benefits from the service provision of a competitor. In these circumstances, a firm will think twice before investing in services because the competitor would have an incentive to avoid the cost of this effort, free ride on the provision of services and offer a lower price. Many authors have shown that RPM can be used to correct for service externalities ([Mathewson and Winter, 1984](#); [Perry and Porter, 1986](#); [Winter, 1993](#)).

Besides, [Dearnley and Feather \(2002\)](#) as well as [Davies et al. \(2004\)](#) find that there is a larger number of brick-and-mortar (B&M) stores in regimes with RPM compared to regimes with free prices. However, [Rey and Stiglitz \(1988, 1994\)](#) point out that vertical restraints that eliminate intra-brand competition can also be used to mitigate inter-brand competition and then would be anti-competitive.

Finally, our article also contributes to the newer literature on machine learning (ML) and text mining approaches. [Varian \(2014\)](#) and [Athey and Imbens \(2019\)](#) provide an overview of important ML methods. [Wang et al. \(2019\)](#) use the *Learning to Place* ML approach to predict book sales and find that a strong driving factor of book sales across all genres is the publishing house. We will use DML techniques similar to the approach of [Knaus \(2021\)](#), who has used DML in the case of musical practice of children on cognitive skills and school performance. For a broad overview on text mining approaches

see [Gentzkow et al. \(2019\)](#). In our article, we will use a latent Dirichlet allocation (LDA) model to determine book genres by analyzing the book descriptions and expert reviews (e.g., [Larsen and Thorsrud, 2019](#)).

3. Data

Our dataset contains prices of e-books for a large number of book titles. In this section, we will present this dataset for our empirical analysis. We first describe the construction of our dataset in Section 3.1 for which we derive descriptive statistics including information on prices, ratings, reviews and the digital size of e-books (in KB) in Section 3.2. In Section 3.3, we present our LDA text mining approach to derive book genres.

3.1. Data Set Construction

The data generating process is structured as follows. We have scraped the *Amazon.co.uk* webpage for a list of publisher and imprint names starting mid February 2020 taking two weeks to get e-book prices as well as further book characteristics available on the Amazon website. Hence, we use the methods of web scraping to generate a cross-sectional dataset. For creating this dataset, we use *a priori* a list of publishing houses, publishers and imprints, which is taken from a historical *Sunday Times* bestseller list. This procedure ensures that our sample only contains books from publishers with a relatively high market share.⁷

⁷The used bestseller list contains entries from January, 2006 until the end of March, 2019.

This proceeding also incorporates books into our dataset which have been published before 2019 because we have done the publisher search on *Amazon.co.uk* independently of the format. Thus, it may have happened that for a certain book title, which we have found within our observation period, another format of the same title has already been published a few years ago. However, we have ensured that no book is included in our working data set which has been published earlier than 2010.

Our raw dataset consists of roughly one million observations, whereby one observation contains several information on different prices, formats, descriptions, ratings, reviews etc. being available on the Amazon website. For every book title there are three entries if all formats (hardcover, paperback, e-book) are available for a certain book title. However, due to the usage of web scraping methods the dataset contains of some entries that are duplicates or not of interest for our analysis. Hence, after the data cleansing process our working dataset consists of 77,629 observations, respectively 47,161 unique book titles.⁸ For our empirical analysis, we drop the hardcover and paperback book titles (see Section 4). Moreover, we only use e-book titles in our estimation approach for which all explanatory variables (book characteristics) are available. Thus, for our empirical analysis 12,001 e-book titles remain in the final working dataset.

Our variables of interest are the retail price, which is the price a consumer must pay for a certain e-book, and the treatment variable *Agency*,

⁸In the former number all three format types (hardcover, paperback, and e-book) are included.

which takes the value one if there is a text field on the Amazon webpage of a book title expressing *'This price was set by the publisher'* and zero otherwise.⁹ Beyond, we have data on several control variables for our empirical analysis. These variables comprise book characteristics as the book format, the book genre, the size of an e-book in KB, variables on book reviews as the star rating and the number of consumer and expert reviews, variables containing information on the author or publisher of a book title and other variables as the publication date or the recommended retail price (RRP). Table 1 summarizes the descriptions for all variables included in our dataset.

Variables	Information
Price	Retail price from the upper right <i>Buy-Box</i>
Format	Hardcover, paperback, Kindle
Star rating	Average rating normalized to be between 0 and 1
No. customer reviews	Number of consumer reviews
No. expert reviews	Number of expert reviews on Amazon
Series	Dummy variable whether book is part of a series
Description and reviews	Detailed text-information on the book and by different reviewers
Genre	Constructed by LDA from the descriptions and reviews (see Section 3.3)
RRP	Recommended retail price which is the print RRP. For Kindle it is either related to the hardcover or paperback RRP
Agency	Dummy variable to be one if the price was set by the publisher and zero otherwise. Only possible for e-books
Seller	Sample is restricted to be sold by Amazon
Author	Information on the author of a book
Title	Information on the title of a book
Kindle.Size	Kindle file size (in KB)
Publisher	Name of the publisher. We have different levels of aggregation (Imprint,Publisher,Publishing House)
Amazon rank	Uncategorized Amazon bestseller rank for either print books or e-books
Bestsellers	Number of bestsellers in the Sunday Times Bestseller List conditional on the Author's name
WeekInChart	Average number of weeks in the bestseller charts conditional on the Author's name
Identifier	Aggregation of ASINs to verify the books
Publication Date	Publication Date of a book title

Table 1: Relevant Variables per book title and the information content they provide.

We have also matched the dataset obtained from Amazon with a historical *Sunday Times* bestseller list to identify authors who have already written a bestselling book title in the past. This variable is important for our empirical analysis (in which we estimate the retail price on an e-book) since the name of a bestseller author is an important quality signal for the

⁹See Figure A.8 in Appendix A for an example from the Amazon webpage.

book readers.

Each book is a unique product written by an author and mostly published by one publisher. Thus, books are heterogeneous goods which make it impossible to actually compare the value of one specific book with one another. In order to provide an acceptable analysis, it is therefore also necessary to control for the genres of the several books. Hence, we use a Latent Dirichlet Allocation (LDA) to derive book genres from the descriptions and reviews of the individual books available on the Amazon webpage. This control variable should be able to capture specific effects between the individual genres. In Section 3.3, this text mining approach will be explained in more detail.

3.2. Descriptive Statistics

Our final sample consists of 47,161 book titles that have been published on *Amazon.co.uk* by the publishers Bloomsbury, Faber, Hachette, HarperCollins, Oxford, Pan Macmillan, Penguin Random House, Scholastic, Simon & Schuster, and a group of smaller publishers between 2010 and 2020. However, in overall there are 77,629 observations in our dataset since there are several formats available for some book titles. Even though the focus of our empirical analysis is on the price of e-books, in this section we also present some descriptive statistics for the book formats hardcover and paperback to show the relationship between those three book formats.

Table 2 offers descriptive statistics of the variables we use for our empirical analysis, summarized by publishers. In addition to e-book retail prices

	Publisher	Bloomsbury	Faber	Hachette	HarperCollins	Small Pub.	Oxford	Pan Macmillan	Penguin Random House	Scholastic	Simon & Schuster
Retail Price	mean	10.37	5.84	5.31	5.42	5.80	16.71	7.79	6.56	5.01	7.45
	std	7.02	3.21	2.86	3.20	6.07	10.12	3.86	2.70	1.58	3.27
Sales Rank	mean	554,409.87	326,982.37	262,402.04	427,905.75	778,429.39	937,724.43	497,575.36	365,604.06	597,692.75	555,114.93
	std	594,258.21	432,500.62	383,524.36	514,257.98	728,316.54	653,612.74	551,601.82	533,573.80	655,036.93	649,582.93
Star Rating	mean	0.90	0.85	0.89	0.89	0.86	0.90	0.89	0.88	0.93	0.90
	std	0.10	0.13	0.08	0.09	0.12	0.10	0.10	0.08	0.07	0.08
No. Customer Reviews	mean	51.10	60.80	119.85	106.88	80.50	17.19	124.70	134.59	76.33	142.89
	std	110.04	113.83	178.01	168.94	161.48	42.95	210.71	195.05	138.88	210.61
Pages	mean	292.57	280.60	388.78	342.59	286.86	412.85	336.34	319.01	231.84	367.91
	std	124.35	191.93	1949.17	330.58	212.83	207.28	121.63	140.63	113.11	256.70
Kindle Size	mean	13,078.55	2,849.00	14,134.92	8,288.36	8,626.52	9,995.83	12,118.73	20,466.37	35,684.97	15,118.97
	std	27,346.77	10,529.46	46,813.98	28,600.24	30,096.99	16,416.27	37,011.08	48,498.42	42,097.93	27,885.95
RRP	mean	15.50	10.28	13.07	12.04	12.08	35.29	13.61	14.59	8.31	13.85
	std	9.76	4.88	5.32	5.90	10.22	30.15	5.68	5.91	3.15	5.29
Date Retail	mean	1.64	2.41	1.81	2.09	1.93	2.97	1.66	2.18	2.11	2.13
	std	1.76	2.62	2.06	2.36	1.91	1.81	1.78	2.25	2.11	2.29
No. Expert Reviews	mean	1.74	2.43	2.33	1.25	1.34	1.14	2.24	1.71	0.65	1.35
	std	0.96	1.24	1.06	1.12	1.51	0.77	1.57	1.40	0.85	1.04

Table 2: Summary Statistics

and the RRP, we also observe several characteristics for each book title such as the sales rank of a title on Amazon, the customer ratings, the number of customer and expert reviews, and the number of pages. As shown in the table, e-books from Scholastic exhibit the lowest average retail price, while the e-books from Bloomsbury have the highest mean prices. Beyond, the titles from Hachette have the lowest average book rank and the book titles published by Simon & Schuster exhibit the highest average number of customer reviews. Most of the other book characteristics are very similar across publishers.

Figure 1 represents the frequency distribution of the retail prices for e-books (top), paperbacks (centre) and hardcover books (bottom) below £100. It is obvious that e-book prices are in a range between £0 (Minimum price is £0.25) and £10, paperback prices concentrate mostly in the £10-£20 interval and hardcover prices are slightly more expensive. While the distributions of e-books and paperbacks are more compressed, the hardcover prices exhibit a higher volatility. Finally, all three price distributions have significant mass points at candidate focal points (e.g., £0.49 (e-books), £9.99 (paperback) and £15.99 (hardcover)).

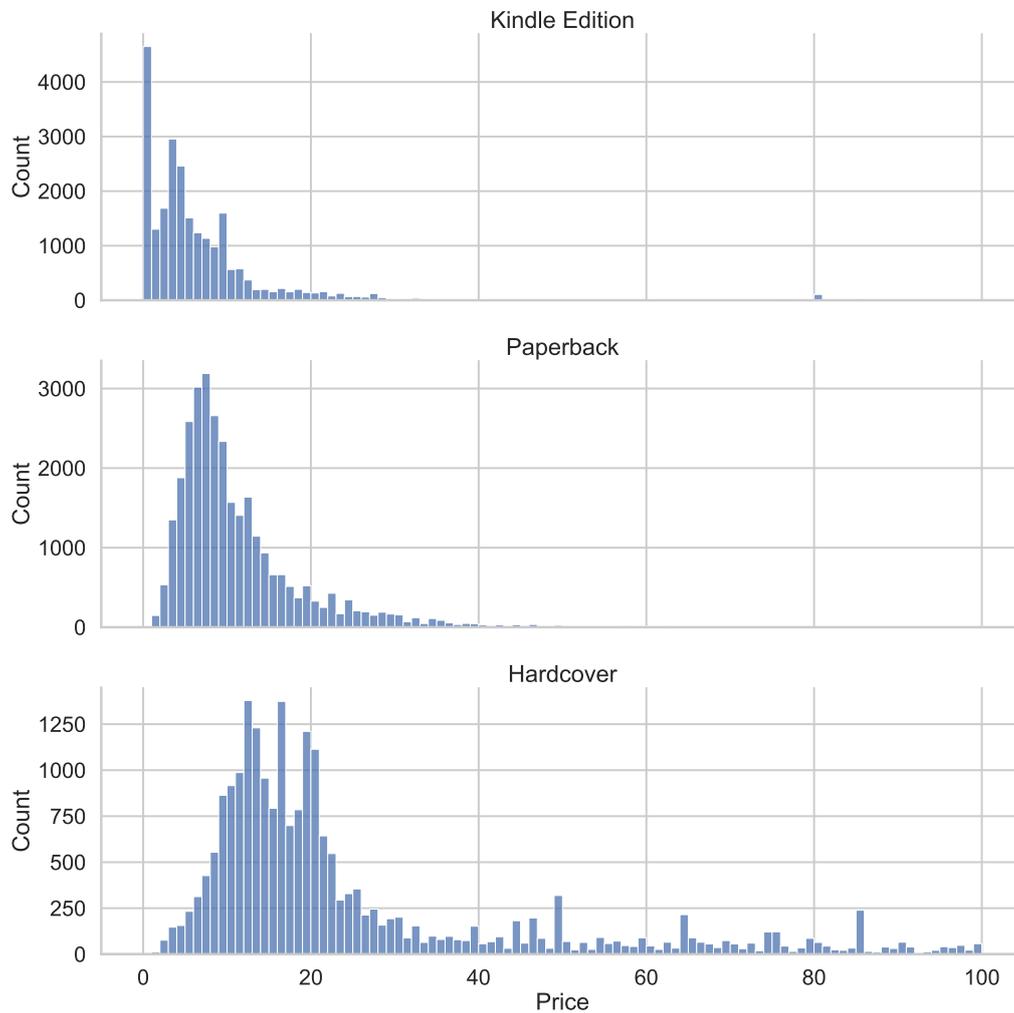


Figure 1: Distribution of retail prices by book format.

Table 3 presents the descriptive statistics for Figure 1. As expected, with an average price of £8.34, e-books are the cheapest of the three book formats, followed by paperbacks (£12.34) and hardcover books (£28.54). The high standard deviation for hardcover books confirms its high volatility, which we have already detected in Figure 1. In overall, the descriptive statistics on book formats imply that hardcover books exhibit the highest

quality of the three formats and confirm the results from [Li \(2019\)](#), who has found that e-books and paperbacks are closer substitutes than e-books and hardcover books.

Format	count	mean	std	min	25%	50%	75%	max
Hardcover	22,647	28.54	42.17	1.05	12.99	18.54	28.96	1575.0
E-book	23,991	8.34	13.89	0.25	2.19	4.99	9.18	467.8
Paperback	30,870	12.35	18.86	0.31	6.55	9.01	13.95	1899.0

Table 3: Retail prices grouped by book format.

The e-books of the individual publishers are sold under different pricing arrangements. While the major publishing houses all have adopted the agency model, Bloomsbury, Faber, Oxford, Scholastic and the group of smaller publishers still use the wholesale model. Amazon mentions on its product pages whether the respective publisher has set the price of an e-book. The Figures [A.8](#) and [A.9](#) in the [Appendix A](#) present examples of this by showing screenshots for the e-book *Elon Musk: How the Billionaire CEO of SpaceX and Tesla is Shaping our Future* as well as for the e-book *Pulse*. In Figure [A.8](#), it can be seen in the first box on the right hand side of the Amazon webpage that the 'price was set by the publisher' so that this is an example for the usage of the agency model. On the contrary, in Figure [A.9](#) this information is missing, which means that Amazon sets the retail price for this e-book and it represents an example for the wholesale model.

Figure [2](#) visualizes the distribution of e-book retail prices by the different

publishers. Prices are obviously more dispersed for book titles published by Bloomsbury and Oxford, whereas the other major publishers mostly have books in the range up to £20. The group of smaller publishers has a significantly larger fraction of e-books in the cheapest price range (about £0.49). In Section 3.3, we will present this figure combined with assigned book genres.

As already mentioned above, the e-books on *Amazon.co.uk* are sold under different pricing arrangements, mostly depending on the publisher. This is illustrated in Table 4 for all publishers in our sample. If the number in the column 'Agency' takes the value one, the respective book titles are sold under the agency model, otherwise the titles are sold under the wholesale model. The table shows that all e-books published by the large publishing houses Hachette, HarperCollins, Penguin Random House and Simon&Schuster are sold under the agency model on Amazon. Pan Macmillan has some American imprints which still use the wholesale model for their e-books. For the e-books in our sample published by Bloomsbury, Faber, Oxford, Scholastic and most of the smaller publishers the wholesale model is used.

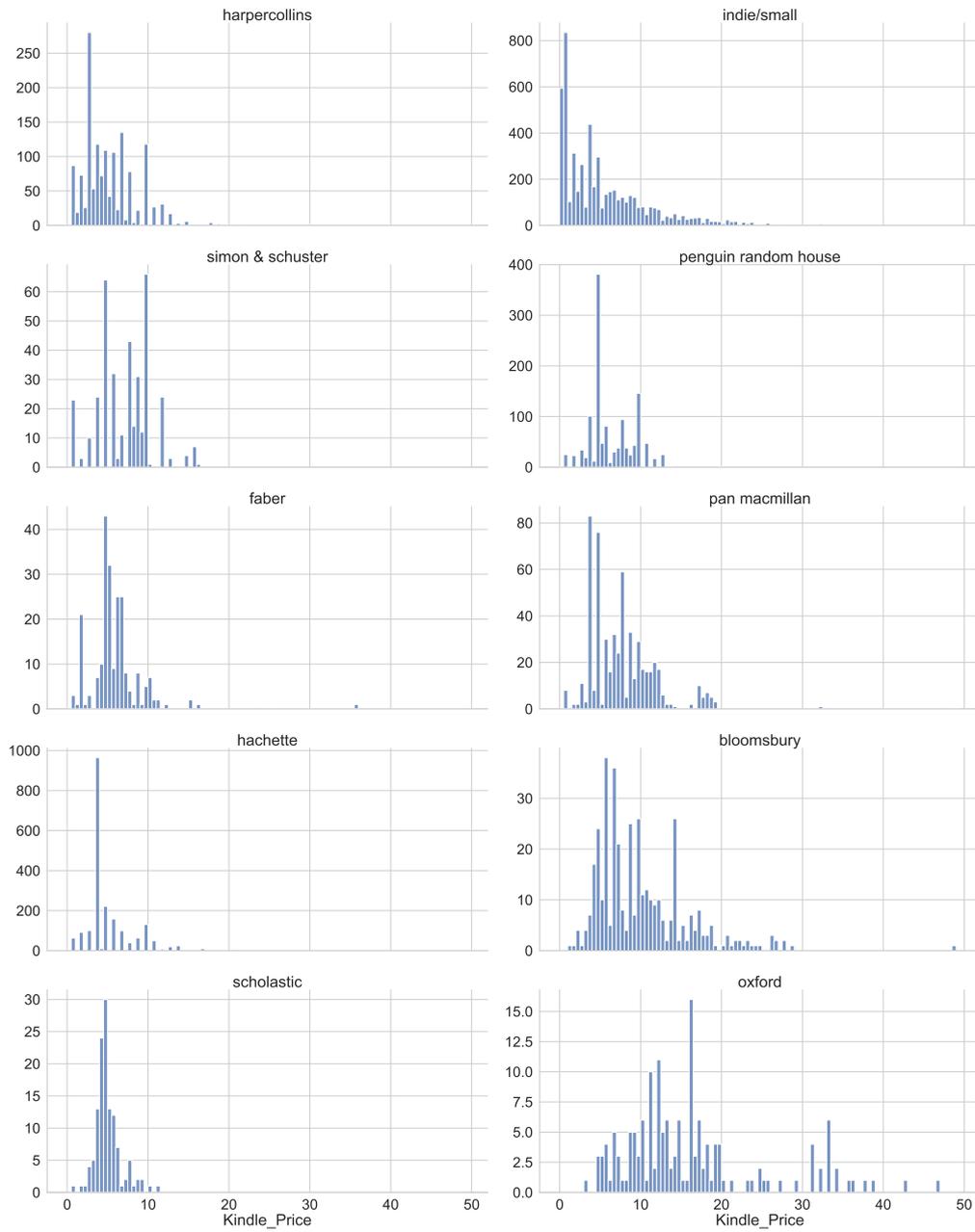


Figure 2: Prices for e-books grouped by publishers. The interval size for each bar is 1 Pound. For illustration purposes the figures are censored at 50 Pound.

Publisher	Agency	Amount	Percentage	Mean Price
Bloomsbury	0	397	100%	10.37
Faber	0	223	100%	5.84
Hachette	1	2,073	100%	5.31
Harper Collins	1	1,469	100%	5.42
Small Pub.	0	4,155	77.37%	6.81
	1	1,215	22.63%	2.32
Oxford	0	161	100%	16.71
Pan Macmillan	0	285	50.35%	9.75
	1	281	49.65%	5.81
Penguin Random House	1	1,240	100%	6.56
Scholastic	0	126	100%	5.01
Simon & Schuster	1	376	100%	7.45

Table 4: Distribution of the agency variable by publishers.

Finally, we want to illustrate the relationship between retail prices for e-books and their book sales rank on Amazon, which is illustrated in Figure 3 by using a scatter plot with a simple regression line. Obviously, there is a positive relation between the retail price and the rank of an e-book in our sample since the regression line has a positive slope. This finding in our sample is in line with the study of Fishwick (2008), who states that 'substantial discounts' (p. 370) have become prominent for bestselling books in the British book market after the abandonment of the Net Book Agreement in 1997.

With respect to the book sales rank, we have to stress how the rank on Amazon is determined. According to sources of Amazon, the ranks are internally updated hourly, but it does not appear immediately. The rank includes current and *all* past sales with higher weights on current sales.¹⁰ This information on the definition of book ranks on Amazon will

¹⁰See https://kdp.amazon.com/en_US/help/topic/G201648140. (Last accessed:

be important for our estimation approach because the price today might be affected by current sales or past sales, but not necessarily vice versa.

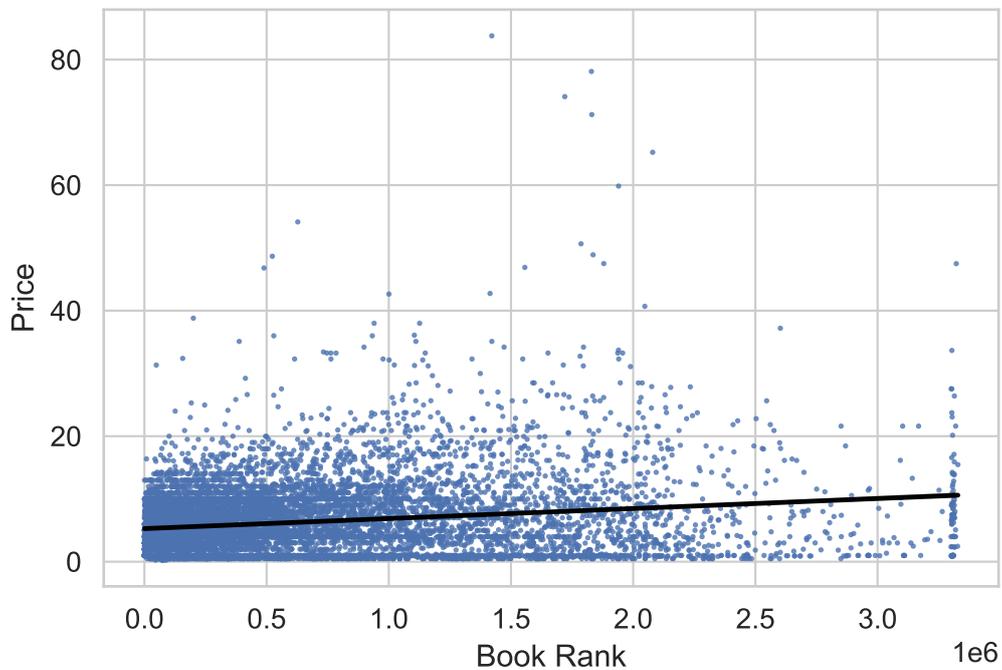


Figure 3: Relation of Amazon book ranks for e-books and retail prices.

3.3. Latent Dirichlet Allocation (LDA)

In the recent past, new technologies have made it possible to use text as data and, therefore, as an input to economic research. Text data, which is inherently high-dimensional, can capture relevant economic concepts not covered by "hard" economic data. In the last years, there has been an explosion of empirical economics research using text as data (e.g., see [Larsen and Thorsrud \(2019\)](#) for an Latent Dirichlet Allocation (LDA) approach or

Lenz and Winker (2020) for paragraph vector topic modelling). We have decided to use an LDA approach to generate book genres and to assign every single book title from our dataset into one of these genres. Such a text mining approach is necessary because on the Amazon webpage the genre information is ambiguous and even not available for some book titles. For this purpose, we use the descriptions and expert reviews from the individual books in our dataset as text data input. We further rely on natural language processing (NLP) to extract the relevant information.

We apply several Python-Modules to clean and prepare the raw dataset.¹¹ Thereby, we remove common words and surnames, eliminate stop words, remove punctuation and pronouns as well as reduce all words to their respective word stems. We note here that around 45,819 unique tokens are kept after this filtering process.

This cleaned descriptions corpus is decomposed into book genres using the already mentioned LDA model. The LDA provides a statistical framework for the generation of documents based on topics. It is an unsupervised topic model that clusters words into topics/ genres, which are distributions over words, while at the same time classifying descriptions as mixtures of topics/ genres. The term "latent" is used because the words are intended to communicate a latent structure, namely, the subject matter (topic) of the description. The term "Dirichlet" is used because the topic mixture is

¹¹Base module is gensim by Řehůřek and Sojka (2010) with a wrapper called Mallet, which is a Java-based open-source NLP text analytics tool (see McCallum (2002) or <http://mallet.cs.umass.edu/>).

drawn from a conjugate Dirichlet prior (Thorsrud, 2020).

The structure of the LDA model is as follows: the whole corpus is represented by M distinct documents (descriptions) and $N = \sum_{m=1}^M N_m$ is the total number of words in all documents. Assuming K latent topics/ genres, each topic is given by a probability vector $\phi_{\mathbf{k}} = (\phi_{k,1}, \dots, \phi_{k,N})$ with $\sum_{n=1}^N \phi_{k,n} = 1$ indicating the probability that each word shows up in this topic. Further, each document $m \in \{1, \dots, M\}$ contains all topics with different probabilities (weights) $\theta_{\mathbf{m}} = (\theta_{m,1}, \dots, \theta_{m,K})$ with $\sum_{k=1}^K \theta_{m,k} = 1$. Both $\phi_{\mathbf{k}}$ and $\theta_{\mathbf{m}}$ are assumed to have conjugate Dirichlet distributions with hyper parameters (vectors) α and β , respectively.

Given $\phi_{\mathbf{k}}$ and $\theta_{\mathbf{m}}$, a document is generated by drawing for each word a topic $k \in \{1, \dots, K\}$ according to the probabilities $\theta_{\mathbf{m}}$ and one word from the selected topic according to its distribution $\phi_{\mathbf{k}}$. This procedure is repeated until the length of the document is reached. To solve the LDA model, we a priori set $\alpha = 50$ and $\beta = 0.01$. The hyper parameter optimization is executed by using Gibbs simulations. Gibbs sampling (also known as alternating conditional sampling) is a specific form of Markov chain Monte Carlo and simulates a high-dimensional distribution by sampling on lower-dimensional subsets of variables where each subset is conditioned on the value of all others (e.g., Steyvers and Griffiths, 2007).

The sampling is done sequentially and proceeds until the sampled values approximate the target distribution. We set the number of sampling iterations equal to 1,000. Then, based on the coherence value across the

estimated LDA models using smaller numbers of genres, we find that 12 topics/genres provide the best statistical decomposition of our book description corpus.¹² A detailed list of all 12 genres is presented in Table B.10 of Appendix B.

One caveat of the LDA estimation procedure is that it does not give the topics/genres any names or labels. Thus, labels are subjectively given to each genre based on the most important words associated with each topic. In the most cases, it is conceptually simple to classify the genres. Besides, the exact labeling plays no material role in our empirical approach, it is just used as a convenient way of referring to the different topics instead of only using topic numbers.

It is more important that the LDA decomposition gives a meaningful and easily interpretable genre classification of the book descriptions, which it does because our LDA approach identifies all important book genres and clearly delineates the topics. This is shown by the Figures 4 and 5, which are two examples of our 12 word clouds to visualize the genre distribution of words by assigned probabilities through the LDA. We have labelled the topic in Figure 4 *Crime Novel/Thriller* and the genre in Figure 5 *Politics*. The larger the size of a word in these clouds is, the higher is its weight within the respective topic.

¹²For 12 different topics, the coherence value exhibits a local peak. We have also considered 9 and 17 different topics, but in the end there was no real change in the effects on the other variables.

visualization purpose in Figure 6, which is based on Figure 2, the largest probability value is chosen to highlight the distribution of topics over e-book prices and publishers. This distribution of prices by genres exhibits the high comparability between the several publishers in our dataset because they are not specialised in certain topics, but all publishers sell book titles from different genres. This is an advantage for our empirical approach because otherwise we would get multicollinearity issues and from an economic point of view we would fail in the sense that we could not compare publishers at all, if there were only specific topics from specific publishers.

Nevertheless, it is obvious that the individual publishers have distinct main topics. For instance, Pan Macmillan primarily publishes fiction titles like crime novels, thrillers or society novels whereas HarperCollins has a focus on the genres family novel and drama. However, it is important not to take these topics at face value because the LDA assigns a probability to each individual topic. That is why we directly include these probabilities as difference to one reference topic in our estimation procedure (see footnote 13 in Section 4).

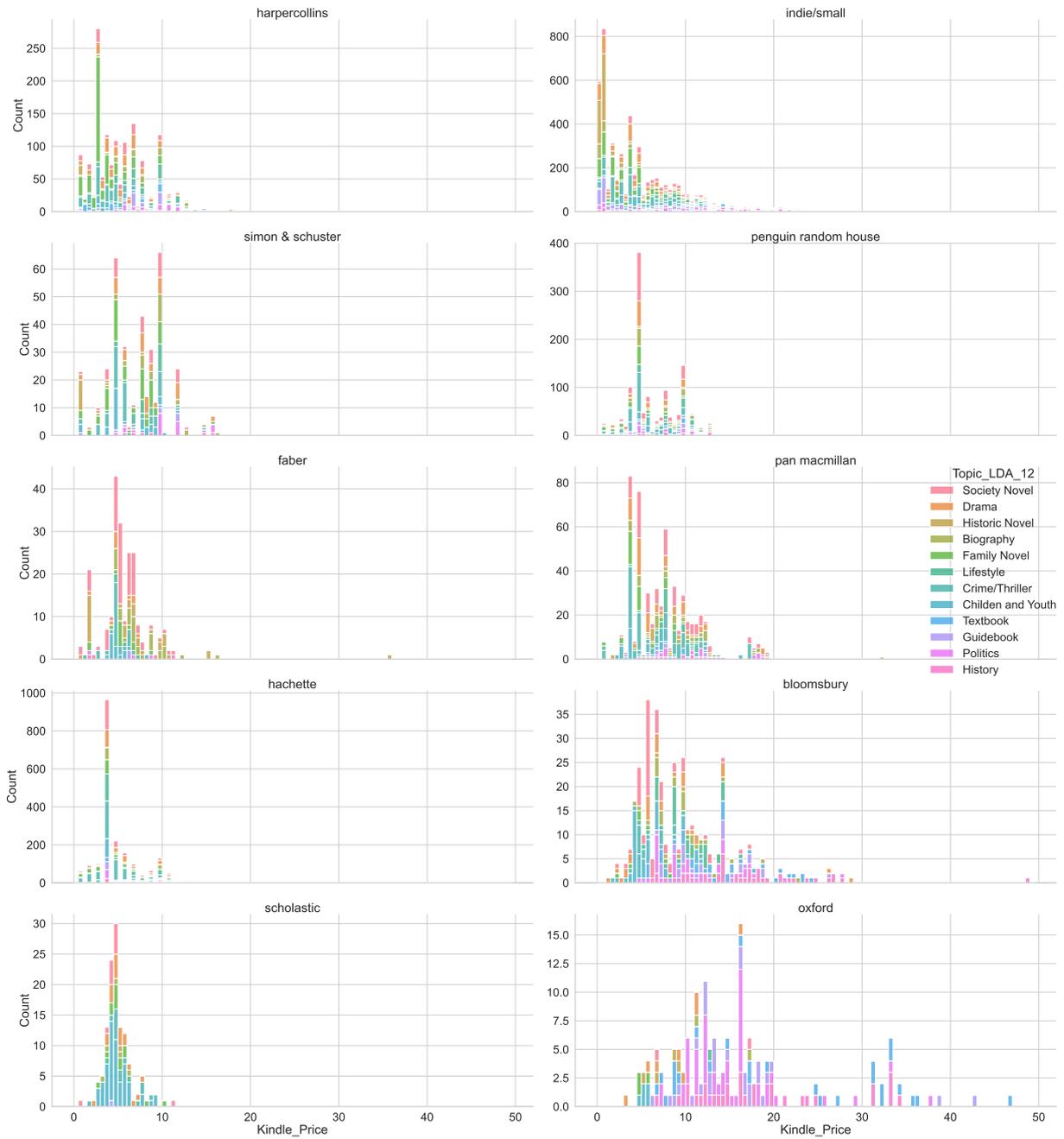


Figure 6: Prices for e-books grouped by publisher and genre. The ordinate is scaled differently for each subplot.

4. Empirical Analysis

In this chapter, we present our main empirical analysis. We first describe our estimation strategy in Section 4.1. In Section 4.2, we present the results of our OLS estimation, in which we estimate the impact of the pricing arrangement on the retail price of an e-book. The propensity score matching is outlined in Section 4.3.

4.1. Estimation Strategy

As already described in Section 1, the goal of our study is to analyze the impact of the pricing arrangement on the retail prices of e-books sold on *Amazon.co.uk*. Therefore, we use publisher and book genre variation of our cross-sectional dataset to estimate the price effect of e-books sold under the agency model. Before turning to the presentation of our estimations, we formalize the hypothesis that is to be tested. If there was no difference between the two types of vertical contracts (agency and wholesale) regarding the price of an e-book, the opportunity for a publisher to set the retail price should (*ceteris paribus*) not have any impact on the final consumer prices. Hence, the hypothesis to be tested is:

Hypothesis 0 (H_0): *The retail price of an e-book is independent of the used vertical contract.*

If there will be a positive correlation between the agency model and the price of e-books, H_0 can be falsified and e-books sold under the agency model are more expensive on average. Observing a negative correlation

would also lead to a falsification of H_0 , but then e-books sold under the agency model would be cheaper on average.

In our baseline estimation approach, we use the standard hedonic modeling approach in the spirit of Rosen (1974), which relies on observing differences in market prices to infer the value or implicit price of underlying characteristics. Thus, we estimate the following log-log OLS model with heteroscedasticity-consistent standard errors:

$$p_i = \alpha_0 + \alpha_1 A_i + \alpha_2 P_i + \alpha_3 G_i + \alpha_4 P_i \times G_i + \alpha_5 R_i + \alpha_6 D_i + \alpha_7 RRP_i + W\theta + \eta_i. \quad (1)$$

In equation (1), the dependent variable p_i is the logarithm of the retail price for an e-book i sold on *Amazon.co.uk*. The treatment variable A_i is a dummy variable which takes the value one if an e-book is sold under the agency model and zero otherwise. P_i contains the publisher fixed-effects, G_i is a continuous variable containing differences of probabilities to a certain reference genre and R_i is a continuous variable for the e-book sales ranks on Amazon. Beyond, D_i reflects the time since a title has been published the first time (in years) and RRP_i gives the recommended retail price of book title i . All other book-specific covariates are included in the matrix W (see Table 1 in Section 3.1).

However, there may be two endogeneity issues in regression equation (1), which might distort our estimated coefficients. First and foremost, our treatment variable A_i might be endogenous due to selection effects.

Specifically, e-books sold under agency contracts might differ from e-books sold under wholesale contracts in ways that were directly related to their retail prices. So it was probably not a random process whether an e-book is sold under the agency or the wholesale model. To eliminate this potential selection bias, we perform a matching procedure in Section 4.3 and use machine learning techniques in Section 5.2.

Second, the rank of an e-book R_i does not only affect the retail price of an e-book p_i but also does the retail price reflect the demand side and, therefore, affect the rank of an e-book which mirrors its sold quantity. As already explained, e-book ranks on *Amazon.co.uk* are internally determined by overall weighted sales ranks. Thus, ranks might be driven by the quantities sold today but the relation is ambivalent since the rank is also affected by quantities sold in the past. The e-book prices can be affected by current and past sales but the impact of prices on total (current and past) sales is not so clear. Nevertheless, we cannot clearly reject the endogeneity issue due to a potential reverse causality. To resolve this potential source of endogeneity between the e-book price and its sales rank, we will also present an instrumental variable approach in Section 5.1.

4.2. OLS Estimation

The results of our main log-log OLS model are outlined in Table 5. We estimate four different specifications for our baseline estimation approach: in the first column of Table 5, we present a naive OLS regression model without publisher and genre effects, in column (II) we include publisher

effects, in the third column we additionally integrate genre effects¹³ and, finally, column (IV) also contains the interaction term $publisher \times genre$.

It is obvious that there is a negative and significant effect of the pricing arrangement *Agency* on the retail price of e-books across all four different specifications. According to the amount, the effect is between 18.5% and 50.3% depending on the exact specification. For the regression in column (I), an e-book which is sold under the agency model on *Amazon.co.uk* is approximately 18% cheaper than an e-book which is sold under the whole-sale model on average.¹⁴ This is the lowest effect across the four different specifications but, however, this coefficient might be biased due to omitted variables because neither publisher nor genre fixed effects are included there, even though these variables are crucial to explain the retail price of an e-book. Including publisher and genre fixed effects as well as their interaction term (see column (IV) in Table 5) increases the agency effect to 45.14% (on amount).

The estimated coefficients for the other controls shown in the table are very similar across the different specifications. The sign of the variable *log sales rank* indicates that e-books with higher sales ranks are sold at higher

¹³ We have described the process to generate book genres by using an LDA approach in Section 3.3. Thereby, we have identified a topic for every e-book title based on the largest probability assigned by the LDA. However, this procedure may be misleading because sometimes the genre probabilities for a book title could be very similar. Therefore, we have solved this issue within our estimation approach by using one reference category and build the difference of the probabilities to the other categories. So by differentiating them, a positive value means that the the respective category has a higher probability assigned by the LDA than the reference category.

¹⁴To calculate the exact effect of the dummy variable *Agency* on the dependent price variable, the formula $100 \times (e^\beta - 1)\%$ must be used.

prices, which we have already demonstrated descriptively in Figure 3 of Section 3.2 and which confirms the results of Fishwick (2008) whereupon bestsellers are sold cheaper in the UK due to 'substantial discounts'. However, we will discuss potential endogeneity issues concerning this control variable in Section 5.1. Moreover, there is a significant and positive relation between the RRP of an e-book and its retail price (see variable *log RRP*), which is not a surprise. We can further observe that the memory space of an e-book (given in KB) has a positive and significant price effect (*log Kindle Size*).

There are five control variables remaining in Table 5. Consumer and expert recommendations seem to drive the price of e-books because the consumer star rating (*log star rating*) as well as the number of expert reviews (*No. expert reviews*) have a positive and significant effect in our regressions. On the contrary, the time since the publication of a book title (in years) seems to be irrelevant for e-book prices since the covariate *Date Retail* is only significant in one of the OLS specifications (column (I) of Table 5).

Lastly, two covariates remain in Table 5 which control for the author's quality. The explanatory variable *WeekInChart* reflects the average number of weeks former bestsellers of an author have last in the bestseller charts of the *Sunday Times* and the continuous variable *Bestsellers* exhibits the number of bestselling book titles an author has written historically. As expected, both variables have a positive and significant effect on the retail price of an e-book in all four specifications, which can be interpreted as

quality signals increasing the price of a book title.¹⁵

	(I)	(II)	(III)	(IV)
Constant	-1.85591*** (0.06562)	-1.51840*** (0.06432)	-1.46891*** (0.06894)	-1.28778*** (0.07080)
Agency	-0.20473*** (0.01124)	-0.69970*** (0.02325)	-0.61211*** (0.02397)	-0.60046*** (0.02461)
log sales rank	0.05257*** (0.00369)	0.07764*** (0.00343)	0.08014*** (0.00364)	0.07919*** (0.00369)
log RRP	1.09110*** (0.01082)	0.95753*** (0.01082)	0.88515*** (0.01187)	0.85294*** (0.01229)
log Kindle Size	0.04495*** (0.00392)	0.03660*** (0.00385)	0.03364*** (0.00423)	0.03617*** (0.00416)
log star rating	0.53052*** (0.04864)	0.38600*** (0.04164)	0.34469*** (0.04112)	0.30686*** (0.03989)
No. expert reviews	0.01513*** (0.00433)	0.01226*** (0.00401)	0.01342*** (0.00399)	0.01192*** (0.00407)
Date Retail	0.01042*** (0.00264)	0.00389 (0.00249)	0.00390 (0.00247)	-0.00102 (0.00248)
WeekInChart	0.00727*** (0.00227)	0.00463** (0.00206)	0.00532*** (0.00206)	0.00491** (0.00197)
Bestsellers	0.00148*** (0.00037)	0.00105*** (0.00035)	0.00115*** (0.00035)	0.00089** (0.00035)
R-squared	0.60288	0.66649	0.67972	0.69700
Adj. R-squared	0.60245	0.66588	0.67884	0.69363
Number of observations	12,001	12,001	12,001	12,001
Publisher	No	Yes	Yes	Yes
Genre	No	No	Yes	Yes
Publisher x Genre	No	No	No	Yes

Table 5: Baseline regressions (log-log OLS). Dependent variable is the logarithm of the retail price for e-books sold on Amazon.

However, the results of our OLS estimation in Table 5 might be biased and inconsistent due to endogeneity issues regarding our treatment variable *Agency* (see Section 4.1 for a discussion). Hence, we will apply a matching procedure in Section 4.3 and use machine learning techniques (see Section 5.2) to reduce the potential selection bias. Beyond, we will follow an in-

¹⁵The variables *WeekInChart* and *Bestsellers* are based on a historical *Sunday Times Bestseller* list. The matching process was conducted via Python’s Fuzzy Matching.

strumental variable approach in Section 5.1 to resolve the potential source of endogeneity of our important control variable book sales rank.

4.3. Matching

A simple OLS estimation would allow us to identify the effect of different pricing arrangements if the treatment was random. Yet, it was probably not a random process whether an e-book is sold under the agency or wholesale model so that there is a selection bias when using a standard OLS estimation model. Hence in a first step, we seek to reduce this selection bias by relying on the propensity score matching (Rubin, 1977; Rosenbaum and Rubin, 1983). Thereby, we identify appropriate treated and control e-books through a matching procedure.

In particular, we identify those e-books that are not sold under the agency model but that had ex-ante the same probability of being sold under the agency model as those that are actually sold under the agency model. To implement the propensity score matching, we first run a logistic regression to recover the likelihood that an e-book is sold under the agency model based on its observable characteristics and use the predicted values from that estimation to collapse those covariates into a single scalar called the propensity score. Second, we match an e-book sold under the wholesale model which is as similar as possible to the considered e-book sold under the agency model based on this propensity score.

The propensity score is the selection probability conditional on the confounding variables ($p(X) = Pr(D = 1|X)$) and relies on two identifying

assumptions. The first assumption is the conditional independence assumption which requires that the outcome variable is independent from the treatment, conditional on the propensity score. This means we should only include variables that are expected to simultaneously influence both the treatment and the outcome. The second assumption is called common support assumption and simply means that for any probability there must be units in both the treatment group and the control group.

Then, the model we estimate through a logit regression as a first step of the matching procedure is:

$$A_i = \alpha + X_i\beta + \eta, \quad (2)$$

where A_i is the probability that an e-book i is sold under the agency model and X_i is a vector of e-book characteristics (see the explanatory variables in equation (1), without publisher fixed effects). In Table 6 we report the estimates for equation (1) on which the propensity scores are estimated. The frequency distributions of the propensity scores for the treated and untreated e-books are presented in Figure 7. As the Figure shows, the frequency distributions of the propensity scores for the two groups of e-books are very similar, indicating that there is a good set of e-books sold under the wholesale model that can be matched with e-books sold under the agency model based on their characteristics. Only a few e-book titles drop out of the sample after the matching process (indicated by the blue and yellow bars in Figure 7).

	Logit
log RRP	-0.0440 (-1.00)
log sales rank	-0.161*** (-11.18)
log star rating	-0.0512 (-0.37)
No. expert reviews	0.344*** (19.24)
log Kindle Size	0.142*** (7.97)
WeekInChart	0.0106 (1.18)
Date Retail	-0.000907 (-0.09)
Bestsellers	0.0193*** (3.57)
Constant	0.691* (2.44)
Observations	12,001
Pseudo R-squared	0.112
Chi-square	1,847.7
Genre	Yes

t statistics in parentheses
* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 6: Logistic regression with *Agency* as the dependent variable.

In a second step, we use the estimated propensity scores to produce a sample where only the matched pairs of e-books remain to estimate the average treatment effect (ATE) of the agency model on the retail price of e-books. The empirical model that we employ is reported in equation (3), where p indicates two paired e-books (given by the matching procedure):

$$\Delta p_p = \alpha_0 + \alpha_1 A + \alpha_2 \Delta X_p + \epsilon_p. \quad (3)$$

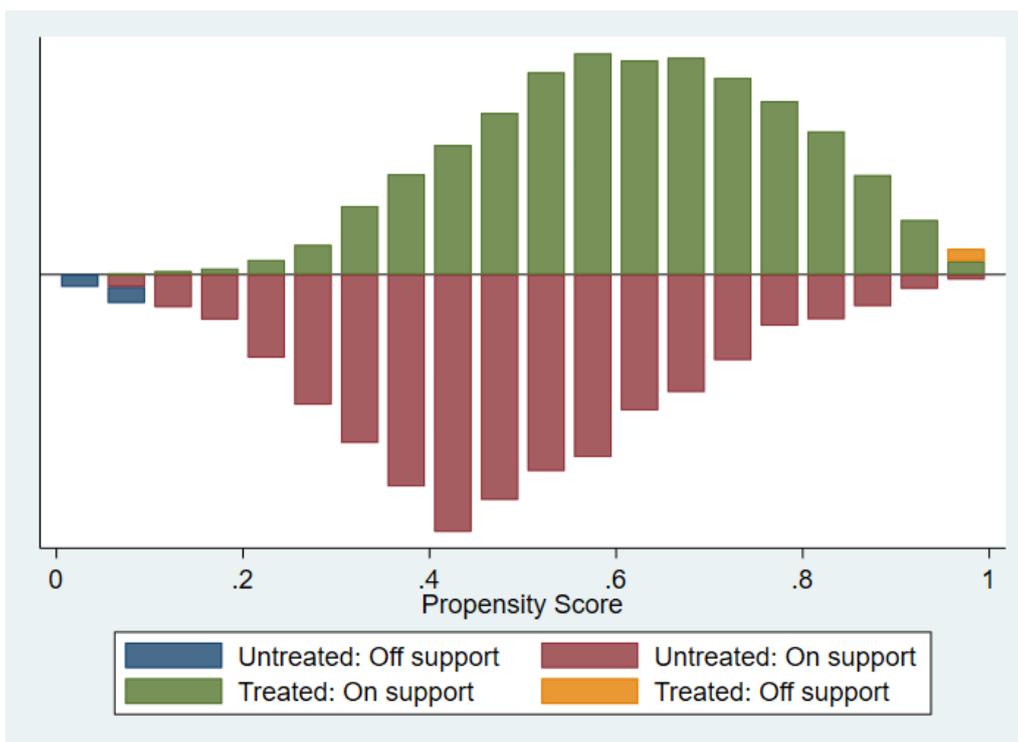


Figure 7: Propensity scores of a logistic model with Gaussian kernel fitting model. We have 11,887 observations which are on support. We present the covariate balance in Table A.10 of Appendix A.

Thus, Δp_p is the difference in the retail price between a treated e-book and a non-treated (control) one. The treatment variable A is equal to one if an e-book is sold under the agency model and zero otherwise. ΔX_p gives the differences in the e-book characteristics between the treatment and the control group.

Table 7 reports the average treatment effect (ATE) from equation (3). We apply four different matching methods (two nearest neighbour covariate

matching procedures, kernel and caliper) to deal with the issue of possible non-randomness in the treatment. The agency model still has a significant and negative effect on the retail price of e-books. The effect ranges between 17.2% and 18% depending on the used matching procedure.

	(1)	(2)	(3)	(4)
	Neighbor5	Neighbor10	Kernel	Caliper
Agency	-0.193*** (-12.82)	-0.192*** (-13.98)	-0.199*** (-17.56)	-0.189*** (-15.47)
OnSupport	11,887	11,887	11,887	11,887

t statistics in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 7: Propensity Score Matching. Standard Errors calculated by bootstrap with 101 random draws. The different models Neighbor5, Neighbor10, Kernel and Caliper refer to a 5-nearest neighbor, 10-nearest neighbor, a Gaussian kernel and a caliper matching, respectively. The caliper radius is determined by 0.2 times one standard deviation from the propensity scores derived from the logistic regression in Table 6 (see Austin, 2011, for optimal caliper length)

5. Robustness Checks

To check the robustness of our results presented in the previous section, we apply two further estimation approaches. First, we use an instrumental variables approach in Section 5.1 to control for the potential endogeneity issue of the variable book sales rank. Following, we present a double machine learning (DML) approach in Section 5.2.

5.1. IV Estimation

The results of our OLS estimation in Section 4.2 might be biased and inconsistent since the important control variable book sales rank might be endogenous. Hence, we will follow an instrumental variable approach in

this section to resolve this potential source of endogeneity. We use the logarithmized number of customer reviews (*log no. customer reviews*) as an instrument for the book sales rank to avoid inconsistent estimates due to reverse causality.

Our instrumental variable *log no. customer reviews* is highly correlated with our endogenous regressor book sales rank (relevance condition) but should have no partial effect on the price of an e-book (orthogonality assumption). Customer reviews can enhance the awareness and information quality for a consumer and, thus, change the tendency for a consumer to purchase a book. However, the absolute number of customer reviews does not affect the purchasing decision of a consumer for a book title, but only surprisingly positive (negative) reviews can increase (decrease) the consumption of a given good (Reimers and Waldfogel, 2021). Hence, the absolute number of customer reviews should also have no partial effect on e-book prices, even though our instrument is highly correlated with the book sales rank (as it is an indicator for past sales).

Following the approach explained above, the linear projection in the first stage regression of our 2SLS estimation can be formalized as follows:

$$R_i = \beta_0 + \beta_1 A_i + \beta_2 P_i + \beta_3 G_i + \beta_4 P_i \times G_i + \beta_5 RRP_i + \beta_6 CR_i + W\theta + \xi_i. \quad (4)$$

In equation (4), the dependent variable R_i refers to the sales rank on *Amazon.co.uk* of book title i . The covariates A_i , P_i , G_i , RRP_i , and W have already been described in the context of our baseline estimation in equation

(1). Our instrumental variable *log no. customer reviews* is displayed by CR_i .

The structural equation of our basic model then takes the following form:

$$p_i = \gamma_0 + \gamma_1 A_i + \gamma_2 P_i + \gamma_3 G_i + \gamma_4 P_i \times G_i + \gamma_5 RRP_i + \gamma_6 \hat{R}_i + W\theta + \varepsilon_i, \quad (5)$$

where the dependent variable p_i is the retail price of e-book i and the fitted values from the first-stage are captured by \hat{R}_i .

The regression results based on equation (5) are presented in Table 8. In the first two columns two classical IV regressions are depicted, whereas for the columns (3)-(5) a Lewbel approach has been applied. The Lewbel approach is a relatively new method, which tackles the issue of endogeneity in linear systems (Lewbel, 2012). This approach serves to identify structural parameters in regression models with endogenous or mismeasured regressors in the absence of traditional identifying information (e.g., external instruments). Thereby, identification is reached by having regressors that are uncorrelated with the product of heteroskedastic errors.

	(1)	(2)	(3)	(4)	(5)
	Price	Price	Price	Price	Price
log sales rank	0.101*** (0.00790)	0.0967*** (0.00733)	0.106*** (0.00884)	0.104*** (0.00623)	0.0937*** (0.00557)
log RRP	0.971*** (0.0126)	0.851*** (0.0123)	0.972*** (0.0129)	0.971*** (0.0126)	0.851*** (0.0123)
log star rating	0.454*** (0.0476)	0.316*** (0.0399)	0.453*** (0.0477)	0.455*** (0.0475)	0.315*** (0.0397)
Date Retail	0.0110*** (0.00258)	-0.00136 (0.00247)	0.0107*** (0.00265)	0.0109*** (0.00258)	-0.00130 (0.00247)
No. expert reviews	0.0200*** (0.00433)	0.0147*** (0.00415)	0.0200*** (0.00431)	0.0204*** (0.00427)	0.0142*** (0.00407)
log Kindle Size	0.0404*** (0.00427)	0.0359*** (0.00414)	0.0401*** (0.00435)	0.0403*** (0.00427)	0.0360*** (0.00414)
Agency	-0.182*** (0.0114)	-0.599*** (0.0245)	-0.181*** (0.0112)	-0.182*** (0.0112)	-0.599*** (0.0245)
Bestsellers	0.00176*** (0.000379)	0.00100** (0.000357)	0.00178*** (0.000378)	0.00178*** (0.000378)	0.000985** (0.000353)
WeekInChart	0.00991*** (0.00245)	0.00603** (0.00208)	0.0102*** (0.00249)	0.0101*** (0.00248)	0.00584** (0.00206)
log no. customer reviews			0.00230 (0.00593)		
Constant	-2.294*** (0.115)	-1.508*** (0.107)	-2.360*** (0.136)	-2.327*** (0.0954)	-1.470*** (0.0889)
Publisher	No	Yes	No	No	Yes
Genre	Yes	Yes	Yes	Yes	Yes
Publisher x Genre	No	Yes	No	No	Yes
Instruments	Reviews	Reviews	Generated	Reviews and Generated	Reviews and Generated
J-Stat			186.6	165.2	310.2
Degree of Freedom for J-Stat			23	23	131
p Value of J-Stat			1.41e-27	1.73e-23	1.45e-16
F-Stat	810.1	242.2	781.9	816.5	243.8
C-Stat	19.44	7.158	16.10	41.49	11.95

Standard errors in parentheses
* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 8: Model Comparison of IV- (1-2) and Lewbel-approaches (3-5). Generated instruments refers to the generation of instruments by the established approach of (Lewbel, 2012).

The results of our IV approach confirm that e-books sold under the agency model on *Amazon.co.uk* are on average significantly cheaper than book titles sold under the wholesale model. Compared to our OLS estimation results in Section 4.2 the estimated coefficients for the variable *Agency* only differ in their magnitude. Also the effects of the other explanatory

variables barely differ between the OLS and the IV estimation approaches, even though the impact of the book sales rank has become larger in the IV regressions.

5.2. Double Machine Learning (DML)

There are further methods beyond the established approaches in the standard econometric analysis for causal inference. We have already applied standard econometric methods as the OLS estimation, the IV approach and the propensity score matching to deal with econometric issues. Recent advances in machine learning approaches also offer a larger toolbox for empirical analyses in economics (e.g., [Athey and Imbens, 2019](#); [Athey, 2018](#), for a broad overview).

Due to recent developments in the machine learning literature there are many approaches, e.g. the DML technique, that gives the possibility to deal with common econometric issues as confounding variables or variable selection by using cross-validation and non-parametric models. This technique also allows to make use of non-standard modelling for relations between variables, like the independent variables have a specific, often assumed linear or quadratic effect on the dependent variable. It permits to use any arbitrary machine learning technique relying on algorithms to find a fitting model for some chosen score functions like the mean squared error.

The reason for relying on further non-parametric/semi-parametric regressions is to circumvent the imposition of a specific model structure. The algorithm of the machine learning models will choose the best fitting model

under some restrictions or parametrization. We will then use regularized linear regression techniques like the least absolute shrinkage and selection operator (Lasso) or regression trees/forests. These methods help to compare our standard econometric approaches with models that are able to ignore irrelevant variables or include non-linear effects (Athey and Imbens, 2019). Moreover, this methods can help in a similar way like propensity score matching to deal with the underlying selection issue (e.g., Lee et al., 2010; Westreich et al., 2010; Knaus, 2021).

Therefore, we apply DML techniques and compare the results to our previous estimations to provide further robustness checks. From a prediction's perspective, the estimations will be split into three different regression models, which is proposed by Chernozhukov et al. (2017, 2018). The models use the DML framework to deal with high dimensional variables, non-parametric functional forms, or unobserved confounding variables which have to be addressed by, e.g., an instrumental variable approach, based on the *DML-Conditional-Average-Treatment-Effect-Estimator*.¹⁶

The equation system we estimate has the following general form (we drop the index i for each book), which stems from the partial linear regression

¹⁶For the estimation implementation we follow the Python Module *econml* provided by Battocchi et al. (2019) (see also <https://econml.azurewebsites.net/spec/estimation/dml.html#overview-of-formal-methodology>).

model of [Robinson \(1988\)](#):

$$\begin{aligned}
 p &= \theta A + q(W) + \varepsilon \\
 A &= f(W) + \eta \\
 \text{s.t. } \mathbf{E}[\varepsilon|W] &= \mathbf{E}[\eta|W] = \mathbf{E}[\varepsilon \cdot \eta|W] = 0,
 \end{aligned} \tag{6}$$

where p denotes the price of an e-book depending on the agency dummy variable A and function $q(W)$ depending on the covariates W . The agency dummy variable A is explained by a function $f(W)$ of these covariates W (similar to the logit specification in equation (2)). The variables η and ε represent stochastic error terms. The conditional expectation functions then can be solved by some non-parametric regressions:

$$\begin{aligned}
 q(W) &= \mathbf{E}[p|W] \\
 f(W) &= \mathbf{E}[A|W].
 \end{aligned} \tag{7}$$

In a next step, the residuals of the price, \tilde{p} , and the residuals of the agency dummy variable, \tilde{A} , are computed by subtracting the fitted values (given by the regression tasks in (7)) from the actual price p and the actual agency dummy variable A :

$$\begin{aligned}
 \tilde{p} &= p - q(W) \\
 \tilde{A} &= A - f(W)
 \end{aligned} \tag{8}$$

Finally, we use these residuals for the prices \tilde{p} and the agency dummy variable \tilde{A} to estimate a linear treatment effect θ that is unbiased based on the assumptions of [Chernozhukov et al. \(2018\)](#):

$$\tilde{p} = \theta\tilde{A} + \epsilon. \tag{9}$$

In Table 9, the column *Model* represents the applied functional form. The first entry in this column refers to the functional form of computing and predicting \tilde{p} and the second one for classifying \tilde{A} . Therefore, *Lin-Logit* relies on OLS and a logistic regression, *Lin-Lasso* relies on an OLS and a logistic regression including L_1 penalty (called Lasso), *Lasso-ElasticNet* uses a Lasso and a logistic regression with a combination of L_1 and L_2 penalties (*Elastic Net*), *Lasso-RFC* refers to a Lasso regression and a random forest classifier, *RFR-ElasticNet* combines a random forest and an *Elastic Net*, *RFR-RFC* uses a random forest for both stages and XGBoost relates to Extreme Gradient Boost for both stages.

The column *Score* displays the mean squared error of the final stage. In the final stage, a simple linear regression is used to get the conditional average treatment effect. There are many more possible estimation techniques but these are sufficient to highlight the stability of our results. The hyperparameters for each model are chosen from a reasonable set and then we use 3 – 5 cross-fold validation within Python’s Sklearn GridSearch. Besides, we also do another 5-fold splitting in each estimation. The presented results outline the best estimation (lowest score) for each model class.

Model	Agency	Std. Error	p-value	Score	Perc.Change
Lin-Logit	-0.1957	0.0112	0.0000	0.3060	-17.7707
Lin-Lasso	-0.1956	0.0112	0.0000	0.3060	-17.7686
Lasso-ElasticNet	-0.1955	0.0112	0.0000	0.3061	-17.7578
Lasso-RFC	-0.2165	0.0126	0.0000	0.3064	-19.4696
RFR-ElasticNet	-0.1607	0.0094	0.0000	0.2134	-14.8437
RFR-RFC	-0.1800	0.0112	0.0000	0.2135	-16.4753
XGBoost	-0.1389	0.0120	0.0000	0.2313	-12.9706

Table 9: Double Machine Learning Approach. The dependent variable is the e-book retail price and the treatment variable is *Agency*. The price effect of the agency arrangement for the respective model is given in the column *Agency* and represents the ATE. Column *Score* refers to the mean squared error. Note: The data has been mean centered with unit variance as Lasso requires this normalization.

The point estimates of the individual regression models are given in the column *Agency* and the relative percentage changes in relation to the intercept are presented in the column *Perc Change* of Table 9.¹⁷ In overall, the DML techniques confirm the results of our main estimations presented in Section 4 and prove the robustness of our regressions, even if we rely on more flexible methods. For instance, e-books sold under the agency model are 16.5% cheaper (on average) than digital books sold under the wholesale model when using the regression model *RFR-RFC*.

6. Conclusion

In this paper, we provide evidence that e-books sold under the agency model on *Amazon.co.uk* are on average significantly cheaper than e-books sold under the wholesale model. Our results are based on an unique dataset containing many characteristics of an e-book. To measure the relationship

¹⁷Again, one can calculate the exact percentage change by using the formula $100 \times (e^{Agency} - 1)\%$.

between the retail price of an e-book and the used pricing arrangement, we have relied on classical econometric techniques as well as on newer methods as the DML approach. We have found an robust and statistically significant effect that e-books sold under the agency model are roughly 18% cheaper than digital books sold under the wholesale model.

The results of our empirical analysis are in line with many theoretical papers studying the price impact of the agency model. Those theoretical analyses argue that retail prices for e-books sold under the agency model are lower due to a lock-in effect exploited by retailers ([Johnson, 2020](#)), the monopolistic power of retailers over a complementary device as it is the case for Amazon when e-books could only be read on a Kindle device ([Gaudin and White, 2014](#)) or because agency selling is just more efficient than the wholesale model and leads to lower retail prices ([Abhishek et al., 2016](#)). Our results also match to the model-theoretical analysis from [Foros et al. \(2017\)](#) if one assumes that competition should be greater among publishers than among retailers, which most likely is the case due to the quasi-monopolistic power of Amazon.

To the best of our knowledge, this paper is the first empirical analysis estimating the price effect of the agency model for e-books by not only incorporating bestselling titles but also "long tail" book titles. Besides, in contrast to previous empirical analyses regarding the retail price of e-books, we apply an LDA approach to determine book genres. Nevertheless, a limitation of our approach is that we use cross-sectional data instead

of panel data so that we cannot control for any dynamic effects on the retail price of e-books. Moreover, we only include one online platform, namely Amazon, instead of comparing various online retailers. Even though Amazon has a relatively high market share for e-books in the UK (see footnote 5), competition between the individual retailers very likely also has an effect on e-book prices.

The dynamic impact of the agency model on e-book prices including bestselling as well as "long tail" book titles remains an open question. Future research should concentrate on panel data to address such dynamic effects of different vertical contracts. Thereby, also other online platforms selling e-books should be included in such an analysis to identify effects between the online retailers. Finally, the long-run impact of the agency model on consumer welfare is an interesting research area. Consumer welfare does not only depend on the e-book prices, but also on other factors such as the number, variety, and quality of book titles written and published.

Acknowledgements

We would like to thank the participants of the following conferences: MaCCI Annual Conference 2021, EARIE Annual Conference 2021 and CRESSE Conference 2021. In particular, we thank Georg Götz, Daniel Herold, Jan Thomas Schäfer, Jona Stinner, Xiang Hui, Franco Mariuzzo and Matthew Olczak for helpful comments. The authors alone are responsible for the content.

References

- Abhishek, V., Jerath, K., Zhang, Z.J., 2016. Agency selling or reselling? channel structures in electronic retailing. *Management Science* 62, 2259–2280.
- Aguiar, L., Waldfogel, J., 2018. Quality predictability and the welfare benefits from new products: Evidence from the digitization of recorded music. *Journal of Political Economy* 126, 492–524.
- Athey, S., 2018. The impact of machine learning on economics, in: *The economics of artificial intelligence: An agenda*. University of Chicago Press, pp. 507–547.
- Athey, S., Imbens, G.W., 2019. Machine learning methods that economists should know about. *Annual Review of Economics* 11, 685–725.
- Austin, P.C., 2011. Optimal caliper widths for propensity-score matching when estimating differences in means and differences in proportions in observational studies. *Pharmaceutical statistics* 10, 150–161.
- Battocchi, K., Dillon, E., Hei, M., Lewis, G., Oka, P., Oprescu, M., Syrgkanis, V., 2019. EconML: A Python Package for ML-Based Heterogeneous Treatment Effects Estimation. URL: <https://github.com/microsoft/EconML>. version 0.12.0.
- Brynjolfsson, E., Hu, Y., Smith, M.D., 2003. Consumer surplus in the digital economy: Estimating the value of increased product variety at online booksellers. *Management science* 49, 1580–1596.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., Robins, J., 2018. Double/debiased machine learning for treatment and structural parameters: Double/debiased machine learning. *The Econometrics Journal* 21.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., Robins, J., et al., 2017. Double/debiased machine learning for treatment and causal parameters. Technical Report.
- Condorelli, D., Galeotti, A., Skreta, V., 2018. Selling through referrals. *Journal of economics & management strategy* 27, 669–685.
- Davies, S., Coles, H., Olczak, M., Pike, C., Wilson, C., 2004. Benefits from competition:

Some illustrative uk cases .

- Dearnley, J., Feather, J., 2002. The uk bookselling trade without resale price maintenance an overview of change 1995–2001. *Publishing research quarterly* 17, 16–31.
- Fishwick, F., 2008. Book prices in the uk since the end of resale price maintenance. *International journal of the economics of business* 15, 359–377.
- Foros, Ø., Kind, H.J., Shaffer, G., 2017. Apple’s agency model and the role of most-favored-nation clauses. *The RAND Journal of Economics* 48, 673–703.
- Gaudin, G., White, A., 2014. On the antitrust economics of the electronic books industry. Available at SSRN 2352495 .
- Gentzkow, M., Kelly, B., Taddy, M., 2019. Text as data. *Journal of Economic Literature* 57, 535–74.
- Gilbert, R.J., 2015. E-books: A tale of digital disruption. *Journal of Economic Perspectives* 29, 165–184.
- Ippolito, P.M., 1991. Resale price maintenance: Empirical evidence from litigation. *The journal of law and economics* 34, 263–294.
- Johnson, J.P., 2017. The agency model and mfn clauses. *The Review of Economic Studies* 84, 1151–1185.
- Johnson, J.P., 2020. The agency and wholesale models in electronic content markets. *International Journal of Industrial Organization* 69, 102581.
- Knaus, M.C., 2021. A double machine learning approach to estimate the effects of musical practice on student’s skills. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 184, 282–300.
- Larsen, V.H., Thorsrud, L.A., 2019. The value of news for economic developments. *Journal of Econometrics* 210, 203 – 218. URL: <http://www.sciencedirect.com/science/article/pii/S0304407618302148>, doi:<https://doi.org/10.1016/j.jeconom.2018.11.013>.
- Lee, B.K., Lessler, J., Stuart, E.A., 2010. Improving propensity score weighting using machine learning. *Statistics in medicine* 29, 337–346.

- Lenz, D., Winker, P., 2020. Measuring the diffusion of innovations with paragraph vector topic models. *PLOS ONE* 15, 1–18. URL: <https://doi.org/10.1371/journal.pone.0226685>, doi:10.1371/journal.pone.0226685.
- Lewbel, A., 2012. Using heteroscedasticity to identify and estimate mismeasured and endogenous regressor models. *Journal of Business & Economic Statistics* 30, 67–80.
- Li, H., 2019. Intertemporal price discrimination with complementary products: E-books and e-readers. *Management Science* 65, 2665–2694.
- MacKay, A., Smith, D.A., 2017. Challenges for empirical research on rpm. *Review of Industrial Organization* 50, 209–220.
- Mathewson, G.F., Winter, R.A., 1984. An economic theory of vertical restraints. *The RAND Journal of Economics* , 27–38.
- McCallum, A.K., 2002. Mallet: A machine learning for language toolkit. <http://mallet.cs.umass.edu> .
- Nelson, P., 1970. Information and consumer behavior. *Journal of political economy* 78, 311–329.
- Perry, M.K., Porter, R.H., 1986. Resale Price Maintenance and Exclusive Territories in the Presence [of] Retail Service Externalities. Department of Economics, State University of New York at Stony Brook.
- Poort, J., van Eijk, N., 2017. Digital fixation: the law and economics of a fixed e-book price. *International Journal of Cultural Policy* 23, 464–481.
- Řehůřek, R., Sojka, P., 2010. Software Framework for Topic Modelling with Large Corpora, in: *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, ELRA, Valletta, Malta. pp. 45–50.
- Reimers, I., Waldfogel, J., 2021. Digitization and pre-purchase information: the causal and welfare impacts of reviews and crowd ratings. *American Economic Review* 111, 1944–71.
- Rey, P., Stiglitz, J., 1988. Vertical restraints and producers' competition. Technical Report. National Bureau of Economic Research.

- Rey, P., Stiglitz, J., 1994. The role of exclusive territories in producers' competition. Technical Report. National Bureau of Economic Research.
- Robinson, P.M., 1988. Root-n-consistent semiparametric regression. *Econometrica: Journal of the Econometric Society* , 931–954.
- Rosen, S., 1974. Hedonic prices and implicit markets: product differentiation in pure competition. *Journal of political economy* 82, 34–55.
- Rosenbaum, P.R., Rubin, D.B., 1983. The central role of the propensity score in observational studies for causal effects. *Biometrika* 70, 41–55.
- Rubin, D.B., 1977. Assignment to treatment group on the basis of a covariate. *Journal of educational Statistics* 2, 1–26.
- De los Santos, B., Wildenbeest, M.R., 2017. E-book pricing and vertical restraints. *Quantitative Marketing and Economics* 15, 85–122.
- Spengler, J.J., 1950. Vertical integration and antitrust policy. *Journal of political economy* 58, 347–352.
- Steyvers, M., Griffiths, T., 2007. Probabilistic topic models. *Handbook of latent semantic analysis* 427, 424–440.
- Telser, L.G., 1960. Why should manufacturers want fair trade? *The journal of law and economics* 3, 86–105.
- Thorsrud, L.A., 2020. Words are the new numbers: A newsy coincident index of the business cycle. *Journal of Business & Economic Statistics* 38, 393–409.
- Tirole, J., 1988. *The theory of industrial organization*. MIT press.
- Varian, H.R., 2014. Big data: New tricks for econometrics. *Journal of Economic Perspectives* 28, 3–28.
- Wang, X., Yucesoy, B., Varol, O., Eliassi-Rad, T., Barabási, A.L., 2019. Success in books: predicting book sales before publication. *EPJ Data Science* 8, 31.
- Westreich, D., Lessler, J., Funk, M.J., 2010. Propensity score estimation: neural networks, support vector machines, decision trees (cart), and meta-classifiers as alternatives to logistic regression. *Journal of clinical epidemiology* 63, 826–833.

- Winter, R.A., 1993. Vertical control and price versus nonprice competition. *The Quarterly Journal of Economics* 108, 61–76.
- Yamey, B.S., 1954. *The economics of resale price maintenance*. Pitman.

Appendix A.

Kindle Books Kindle Unlimited Prime Reading Kindle Book Deals Best Sellers & more Free Reading Apps Buy A Kindle Newstand Audible Audiobooks

Kindle eBooks > Biography & True Accounts > Professionals & Academics

Elon Musk: How the Billionaire CEO of SpaceX and Tesla is Shaping our Future Kindle Edition
by Ashlee Vance (Author) | Format: Kindle Edition
★★★★☆ 6,985 ratings

See all formats and editions

Kindle Edition £5.49 Read with Our Free App	Audiobook £0.00 Free with your Audible trial	Hardcover £13.92 13 Used from £10.12	Paperback £8.19 <i>prime</i> 15 Used from £3.27 22 New from £6.91	Audio CD from £64.95 1 Used from £64.95
--	---	---	---	--

South African born Elon Musk is the renowned entrepreneur and innovator behind PayPal, SpaceX, Tesla, and SolarCity. Musk wants to save our planet; he wants to send citizens into space, to form a colony on Mars; he wants to make money while doing these things; and he wants us all to know about it. He is the real-life inspiration for the Iron Man series of films starring Robert Downey Junior.

The personal tale of Musk's life rimes with all the transmittive associates with a most dramatic flavor

Length: 400 pages | Word Wise: Enabled | Enhanced Typesetting: Enabled | Page Flip: Enabled | Due to its large file size, this book may take longer to download

The 1-Page Marketing Plan: Get New Customers, Make More Money, And Stand Out From The Crowd by Allan Dib
★★★★☆ 1,801 | £1.99 | Kindle Edition
Shop now

Print List Price: £9.99
Kindle Price: **£5.49**
Save £4.50 (45%)
Sold by: Amazon Media EU S à r.l.
This price was set by the publisher.

Buy now

Deliver to your Kindle or other device

Buy for others
Give as a gift or purchase for a group.
Learn more

Quantity: 1 | Continue

Send a free sample
Deliver to your Kindle or other device

Add to List

Enter a promotion code or Gift Card
Share | Embed

READ ON ANY DEVICE

Figure A.8: Screenshot of *Elon Musk: How the Billionaire CEO of SpaceX and Tesla is Shaping our Future* (Amazon.co.uk).

Kindle eBooks > Crime, Thriller & Mystery > Crime Fiction

Pulse Kindle Edition
by Michael Harvey (Author) | Format: Kindle Edition
★★★★☆ 54 ratings

See all formats and editions

Kindle Edition £6.01 Read with Our Free App	Audiobook £0.00 Free with your Audible trial	Hardcover £14.99 4 Used from £4.89 6 New from £2.10	Paperback £5.71 5 Used from £1.97 5 New from £5.71	Audio CD £4.44 1 New from £4.44
--	---	---	--	--

Stranger Things meets The Departed. A hugely original thriller, with film rights already snapped up by New Line Cinema

Daniel Fitzsimmons seems like an ordinary American teenager. He loves his brother, his friends, his school work. But Daniel has powers. Powers that he does not understand and is not sure he can control.

Length: 364 pages | Word Wise: Enabled | Audible Narration: Ready

The Sign of the Blood (A Dangerous Emperor Book 1)
by Laurence O'Brien
★★★★☆ 133 | £0.99 | Kindle Edition
Shop now

Print List Price: £7.99
Kindle Price: **£6.01**
Save £1.98 (25%)

Buy now

Add Audible narration to your purchase for just £4.49
Deliver to your Kindle or other device

Buy for others
Give as a gift or purchase for a group.
Learn more

Quantity: 1 | Continue

Send a free sample
Deliver to your Kindle or other device

Add to List

Enter a promotion code or Gift Card

Figure A.9: Screenshot of *Pulse* (Amazon.co.uk).

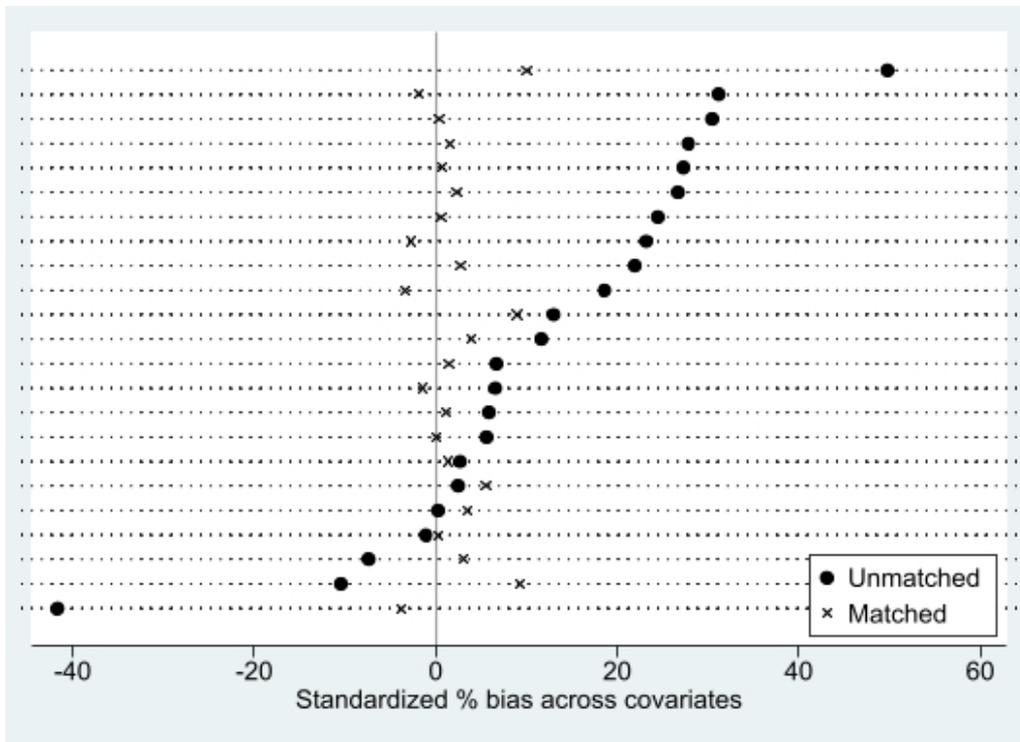


Figure A.10: Covariate Balance. Matched sample exhibits a mean (median) bias of 3.1. (2.3) Rubin's B and Rubin's R are 20.1 and 1.5, respectively.

Appendix B.

Topic	Genre
0	History
1	Guidebook
2	Children and Youth
3	Society Novel
4	Lifestyle
5	Crime Novels/Thriller
6	Politics
7	Historic Novel
8	Drama
9	Family Novel
10	Biography
11	Textbook

Table B.10: 12 different genres identified by our LDA approach.