

A Work Project, presented as part of the requirements for the Award of a Master's degree in **Business Analytics** from the Nova School of Business and Economics.

## **Drivers and prediction of organic search engine CTR**

Improving and understanding predictions through model interpretation

PATRICK GUNDLACH (49594)

Work project carried out under the supervision of:

Qiwei Han, PhD

**25-01-2023**

### **Abstract**

33% of web traffic in the \$5.7 trillion e-commerce industry originates from organic search engine results. Thus, website providers benefit from a holistic understanding of the drivers of click-through rates (CTR) on organic searches to increase traffic. However, providers face a

knowledge gap, as existing literature focuses on position as the primary CTR influence, disregarding other result page characteristics. To solve this problem, I use an extensive dataset comprising organic Google result page information. I conduct an elaborate data analysis highlighting the impact of four categories of result page characteristics before determining suitable CTR prediction modeling techniques. I discover novel patterns impacting CTR for each category and find tree-based models to outperform state-of-the-art deep-learning models. Additionally, by interpreting the XGBoost model, I find potentials for model improvements, quantify the relevance of result page characteristics, and discover further drivers of CTR.

### **Keywords**

organic click-through rate, CTR prediction, e-commerce, SHAP

### **Acknowledgements**

I am very grateful for the continued support of our supervisors, Qiwei Han and Maximilian Kaiser. Additionally, I highly appreciate Grips enabling us to conduct this work by providing the extensive data set used.

This work used infrastructure and resources funded by Fundação para a Ciência e a Tecnologia (UID/ECO/00124/2013, UID/ECO/00124/2019 and Social Sciences DataLab, Project 22209), POR Lisboa (LISBOA-01-0145-FEDER-007722 and Social Sciences DataLab, Project 22209) and POR Norte (Social Sciences DataLab, Project 22209).

# 1. Introduction

Over the past years, the e-commerce business has seen extraordinary growth, with the Covid-19 pandemic additionally boosting its expansion. In 2022, Statista estimates global e-commerce sales at 5.7 trillion US-Dollars (Chevalier 2022a) which exceeds the 2021 GDP of the world's third-largest economy, Japan (O'Neill 2022). With this immense potential, online businesses are flocking the market, trying to gain a share of the steadily growing pie of e-commerce revenue. This revenue is created through sales on a website leading to the ultimate challenge of generating as much traffic on a website as possible. A prominent way to accomplish this is search engine optimization (SEO), an online marketing strategy with the goal of achieving the highest possible traffic for a website (Ledford 2009).

For e-commerce, 33% of the total web traffic is estimated to originate from organic traffic (Chevalier 2022b). Organic traffic refers to cases in which a user searches a keyword in a search engine and arrives at a website by clicking on one of the non-ad results. Because of the high relevance of this organic traffic, it is critical for website providers to improve their chances of generating clicks through search engines. In general, the clicks to a website on a search result page for a specific keyword can be formulated as

$$(1) \text{ Clicks} = \text{Impressions} * \text{CTR}$$

where CTR refers to the click-through-rate of a result (Google 2022a). Accordingly, to achieve higher web traffic, SEO practitioners want to increase at least one of both factors. While the impressions for any keyword can be approximated with the search volume available (Semrush 2022a), the CTR is only available to a website owner for keywords they already rank for. Thus, website providers can easily understand which other keywords are most promising to rank for but usually have very little insight into how much of the volume will

convert into actual traffic to their website. Providing an understanding of what influences CTR can overcome this hurdle and enable companies to receive more website traffic by achieving a higher CTR.

To increase CTR, SEO efforts tend to focus on the position of a website on a result page, as top-ranked websites receive more user attention and clicks (Lewandowski, Sünkler, and Yagci 2020). To rank high on a result page, a website needs to be considered as relevant for the particular keyword. To achieve this, SEO practitioners concentrate their efforts on matching a website's content and metadata to the keyword. (Cui and Hu 2011; Shih et al. 2013; Bala and Verma 2018; Ziakis et al. 2019; Das 2021; Olson et al. 2021)

This approach of focusing only on the ranking of a website on a result page to increase CTR appears one-dimensional, as it ignores how other characteristics of a search page directly influence user behavior leading to the desired clicks.

These characteristics can be grouped into four categories potentially influencing CTR:

1. The *position* of a result on a particular result page.
2. *Keyword* characteristics: These characteristics can relate to a keyword itself (e.g. length or content), the intent of the search, or the accompanying metadata about the result page, such as the search volume or difficulty of ranking high for a keyword.
3. Characteristics of a *result*: Such include URL, title, or description of a result.
4. *SERP features*: These are visual elements containing information from organic search results and aim at making Google result pages more engaging (Google 2022b). These increasingly prominent elements of a search engine result page (SERP) can, for example, take the form of an instant answer to a user's question, a knowledge panel on the right of the result page, or an image next to a result. Appendix I shows an illustration of SERP features.

The categories demonstrate that apart from the position, a website provider can influence various variables to gain users' attention. In order to do so and accurately predict CTR, however, one first needs to understand the characteristics' effects and dynamics. This opens two research questions:

**RQ1:** How do these characteristics influence CTR, and how do they impact the importance of a website's position?

**RQ2:** Which result page characteristics and machine learning models can be used to accurately predict the organic CTR of a website on a Google result page?

In this work, we extend the existing analysis of position as the main direct influence on CTR to additional keyword, result, and SERP feature characteristics summarized in a unique dataset. An analysis including these characteristics provides a novelty in CTR research. In addition to that, we will identify practical CTR prediction model techniques to apply to this data. This enables ex-ante evaluation of the effectiveness of SEO efforts by estimating expected website traffic through predicted CTR and commonly known impressions.

In the following, the second chapter will provide an introductory analysis of related publications on CTR influences and CTR prediction models. The third chapter introduces the used dataset and applied methodology in detail. After presenting the findings of an exploratory data analysis in chapter 4, chapter 5 examines which CTR drivers have the most predictive power and what models perform best for the prediction of CTR.

## **2. Related works**

When raising the question of how to predict CTR and which features impact it, one is mainly looking at two related research fields. First, current research into drivers of organic CTR, and

second, recent research on models and data used to predict CTR. Therefore, this section will provide a high-level literature analysis of the former.

## **2.1 Research into drivers of organic CTR**

Search engine optimization (SEO) is arguably one of the most wide-ranging online advertising strategies today (Olson et al. 2021; Kingsnorth 2022). According to the Oxford dictionary, search engine optimization is defined as ‘the process of maximizing the number of visitors to a particular website by ensuring that the site appears high on the list of results returned by a search engine’ (Stevenson 2010). The definition stresses the direct relationship between traffic and a website's position on result pages. This importance of position is also expressed by the vast majority of research on SEO marketing. (Cui and Hu 2011; Shih et al. 2013; Bala and Verma 2018; Ziakis et al. 2019; Das 2021; Olson et al. 2021). SEO achieves higher ranks by precisely tailoring a website's content to a given keyword, for example, by adding new and qualitative blog posts with the aim to appear more relevant to search engines (Das 2021). Other aspects such as keyword, result, and SERP feature characteristics, if considered, are only analyzed regarding their influence on the position, but not their direct influence on CTR.

To conclude, previous research sees the position as the primary driver of CTR and the most effective way to gain more user attention. We contribute to this perspective by introducing additional characteristics of result, keyword, and SERP features to the analysis of the directly influencing factors of CTR. In the following, we put a special emphasis on SERP features as they are one of the main visual elements that users perceive on a result page, next to organic and paid search results. We will do so by first outlining how they have been described in related works and to what criticism they are linked in public.

While the literature focuses on a website's rank, SERP features are receiving relatively little attention, with little literature researching isolated SERP features. An example is Sam-Martins (2020) analysis of the *featured snippet*, which in turn relies on information found on blogs from influential SEO companies like Semrush. While there is little research on the impact of SERP features, there are a number of influential marketing blogs that advertise the importance of appearing in SERP features (Moz 2022; Semrush 2022b; Wheelhouse 2022). Most blogs argue that SERP features benefit websites featured in them and harm the CTR of websites that appear next to SERP features without being featured (Moz 2022, Wheelhouse 2022). Google itself advertises significant improvements in click-through rate, visits and time spent on websites when chosen to be shown in SERP features (Google 2022b). Critics, however, state that by adding additional features to its result pages' and leveraging information originally published on third-party websites, Google reduces reasons to leave its ecosystem (Tober 2022). This can harm the publishers of the original information that rely on traffic to their websites. Zero-click searches, which are searches where a user does not leave Google's ecosystem, amount to 25.6% of all searches, according to a non-academic study (Tober 2022). This shows that by analyzing the effect of SERP features on CTR, this work has the potential to shed light on the intransparent role of SERP features in Google's search engine ecosystem and provide first empirical research on the SERP features' impact on website performance.

## **2.2 CTR prediction research**

This chapter discusses achievements in previous research on CTR prediction and points out the difference in approaches between previous publications and this work. These differences lie in the nature of the prediction, i.e. whether it is a classification of clicks from individual users or a regression for aggregated CTR, and in the nature of clicks, i.e. whether they are organic or advertisements.

The problem formulation addressed by the vast majority of CTR research is: “Will user X click item Y?”. More technically, the underlying problem is a binary classification problem, i.e. predicting yes or no. The item in question could be an ad on a search engine result page but also an article in an online shop (Yang and Zhai 2022).

This poses a fundamental difference to the problem formulation of this work, as this work deals with a regression problem. It predicts a continuous value based on aggregated data while existing research conducts a classification for a single observation. This difference is also visible in the data used to make predictions. The dataset provided by Grips has high density and low sparsity. Opposingly, the datasets of the existing research, such as the commonly used Criteo dataset, usually contain the columns *user* and *item* with extremely high cardinality for each (Yang and Zhai 2022; Juan et al. 2016; Zhou et al. 2018). To make these columns usable for machine learning models, the values are usually one-hot-encoded, resulting in extremely sparse and high-dimensional feature vectors. Consequently, models are built in such a way that they can incorporate this high dimensionality and sparsity (Zhou et al. 2018).

Recent literature focuses on creating and exploring models suitable for such characteristics. At the beginning of user-item CTR prediction, the research focused on multivariate statistical models like linear regressions enhanced with polynomial features to include feature interaction (Wang, Suphamitmongkol, and Wang 2013; Yan et al. 2014). As these models show difficulties dealing with sparse data, the research progressed to models more suitable for this kind of challenge with a focus on factorization machines and later field-aware factorization machines (Ma et al. 2016; Pan et al. 2018; Yuchin et al. 2016). While these consistently outperform the multivariate models, the literature points out difficulties with generalization (Yang and Zhai 2022). With the rise of convolutional neural networks (CNN), research increasingly attempts to utilize them for CTR prediction (Chen et al. 2016;



Gharibshah et al. 2020). Based on Google's Wide&Deep recommender system architecture (Cheng et al. 2016), the introduction of DeepFM (Guo et al. 2017) marked a new milestone in user-item CTR prediction. DeepFM combines the previously well-working factorization machines with neural networks and lays the foundation for current research. The most recent literature focuses on more accurately modeling feature interactions and produces promising models, such as MaskNet (Wang, She, and Zhang 2021) and DeepLight (Deng et al. 2021), that continuously, though only slightly, outperform DeepFM. Yang and Zhai (2022) provide a detailed overview of the field of CTR prediction in the context of user-item interactions.

Next to the differences between individual click prediction (classification) and aggregated CTR prediction (regression) as described above, a further differentiation lies in the nature of clicks. There are some papers that focus on search engine CTR (Richardson, Dominowska, and Ragno, 2007; Graepel et al. 2010; Zhang et al. 2021). However, only the work of Richardson, Dominowska, and Ragno (2007) focuses on the prediction of aggregated CTR, while the other works predict individual, non-aggregated user interaction. Furthermore, all works focus on paid instead of organic results. This lack of research on organic CTR can be explained by the absence of publicly available data on this topic, as the data is one of the search engines' most significant assets, which therefore has no interest in sharing it.

To summarize, many papers focus on predicting whether an individual user will click on a specific ad and use this information as a recommender system. Additionally, some papers predict aggregated advertisement CTR. However, to the best of our knowledge, there are no publications that predict aggregated organic search engine CTR.

## **2.3 Implications**

The analysis of the research on SEO marketing shows that the position is the most important driver of CTR. Therefore, it promises to play a critical role in predicting organic CTR. In

contrast, the effect of SERP features is theorized but barely analyzed. Therefore, it is essential first to analyze and fully understand the drivers of CTR and potential modeling approaches before predicting CTR. Therefore, in chapter 4, an in-depth exploratory data analysis was conducted to provide fundamental insights into the challenge of CTR prediction on organic search engine result data.

Additionally, while there is a large amount of literature on CTR prediction regarding user-item interactions, there is little research on aggregated search engine CTR prediction. Therefore, while it is insightful to apply the previously mentioned models to the data at hand, it is crucial to consider that the nature of the problem and the available data are different. This is likely to require distinct modeling approaches, which are tested in chapter 5.

### **3. Dataset & Methodology**

This section provides an introductory overview of the used data and methodology applied to answer the research question.

#### **3.1 Dataset**

The dataset used for the analysis was provided by Grips, formerly Peekd, a German start-up that aims at ‘shining a light on the blind spots [in e-commerce sales] and create the most comprehensive map of online commerce for retailers and brands’ (Grips 2022). The data comprises historic Google search data for a given URL and keyword. One row can, for example, resemble metrics related to the performance of the URL ‘www.novasbe.unl.pt/en/’ for the searched keyword ‘nova university’. As such, the data is provided in a tabular, heterogenous format. It ranges over 43 domains, collected from US-based desktop searches between May 31st, 2022, and August 18th, 2022. In the raw dataset, 70 features are available for 79,853 data points.

Although the dataset was provided by Grips, it was enriched with features from Semrush, a company specializing in search engine marketing, and Google's Keyword Planner (see the source of each feature in appendix II.). The features cover various aspects, ranging from the searched keyword over metadata about the result page to result-specific information and which SERP features are shown. Finally, the data includes the click-through rate for each result which, in line with the research question, is determined as the target value.

The dataset, with its comprehensive set of features consisting of real-world, not publicly accessible data for dozens of e-commerce stores, signifies a novelty in available data.

### 3.2 Feature subsets

For a better understanding of the effect of position, keyword, result, and SERP feature characteristics on CTR, the variables in the dataset were grouped accordingly into four feature subsets. If applicable, subsets were enriched with self-generated features. The generated feature subsets form the basic structure for the analysis in chapter 4 and 5. A data dictionary describing all features can be found in appendix II. and a list of all features for each subset in appendix III.

**Position** - The position subset groups features that are exclusively position related. Those are the position of the result during the measurement, the monthly position average, and the difference in positions to the last measurement. While the position subset is technically also handled as a subset in the following analysis, it is important to note that the feature position plays an overarching role in the problem at hand. The literature review and findings of this work highlight that position is the most important feature for predicting CTR. Therefore, in the analysis, the position subset is often used in conjunction with other subsets.

**Keyword** - The keyword subset consists of the keyword itself and everything directly resulting from the keyword a user types into the search engine. This includes information about the keyword, such as the length of the keyword, information indirectly related to that keyword, like the competition or the search volume, and finally, the search intent. One generally differentiates between four intents: (i) The informational search, where a consumer is looking for information about a subject; (ii) the navigational intent in which a specific website is searched, e.g. a marketplace for a product; (iii) the commercial intent that implies a purchasing decision but it is not yet decided on the final transaction (e.g., “linux alternative”); (iv) the transactional intent where a consumer has an action in mind such as “buy iPhone”.

To quantify the complexity of each keyword, the Flesch reading ease score (Flesch 1948) was computed. Even though newer scores have been developed, Flesch’s reading ease score is used because it is commonly accepted and easily interpretable. It counts the average number of words per sentence and the average number of syllables per word. A high score corresponds to good readability, characterized by short sentences and short words.

**SERP features** - The SERP features in the dataset are binary features stating for each entry whether each SERP feature is present (1) or not (0). They can be divided into page and positional SERP features. Page SERP features describe whether a certain feature is present somewhere on the page. Positional SERP features describe whether a result is featured in a particular SERP feature. This can take different forms. For example, a *Review* can be shown beneath a URL, or a *Snippet* from the URL can be displayed in the *Knowledge Panel*. Therefore, if a positional SERP feature is present, the equivalent page SERP feature must be present, but not vice versa. The page SERP features also include information about the kind of ads present on a result page. We aggregated this information to a binary feature describing whether ads are displayed. To add a layer of interpretability, the dataset was additionally enhanced with the total count of SERP features both for page and positional features.

**Result** - The result subset includes all features related to the characteristics of one URL on a result page. While for the three other subsets the features available in the dataset are very comprehensive, available result-related features only cover the URL of the result but not other relevant aspects such as the title or description. To generate informative value out of the URLs, the associated domain is introduced as a one-hot-encoded feature, and two keyword-related features are computed: It is checked whether at least one word of the keyword appears in (1) the domain and (2) the URL.

### 3.3 Data preparation

Before looking into the data in more detail, it has to be prepared for further analysis. To do so, the Google data preprocessing guidelines were followed (Google 2022c). Duplicate rows and columns as well as columns that had only one unique value or were redundant were dropped. Since  $CTR = \frac{clicks}{impressions}$ , both ‘clicks’ and ‘impressions’ are very closely related to the target variable and are consequently dropped, avoiding collinearity and information spillage. Finally, column data types were changed in order to facilitate effective analysis. To minimize noise through observations of result pages that were shown to too few users, an impression threshold of 20 impressions per result was introduced. 58,898 rows remain for analysis. For the modeling, all non-binary features are standard-scaled to a mean value of 0 and standard deviation of 1.

## 4. Drivers of CTR (EDA)

This chapter aims at identifying drivers of CTR through an exploratory data analysis (EDA). We focus on non-obvious and novel insights as well as characteristics that are relevant for CTR prediction models. Thus, the EDA does not only create a basic understanding of the data but also as the foundation for selecting and examining predictive models in chapter 5.

The analysis is structured along the four feature subsets and dedicates a section to each. In addition to that, we analyze feature interactions and present a summary of the main findings and implications for the modeling.

#### **4.1 Position - the most important predictor**

There are two main findings regarding the *position* subset: First, the position is more relevant than other subsets in determining CTR. Second, results on position one receive much higher average CTRs than results on any other position.

In line with previous research (see chapter 2 and Iqbal et al., 2022, 4), our data shows that of all features, the position has the strongest impact on CTR. This is based on position and related features having the highest correlations with CTR. The average position per month has a Pearson correlation of 0.5 with CTR. Related features, such as the previous position, have a Pearson correlation of more than 0.4. This is significantly higher than the second most correlated variable with CTR ( $\rho = 0.25$ )<sup>1</sup>.

Our findings show that the CTR on position one is much higher than on other positions, as many search engine users only read and click the first result. The result on the first position receives an average CTR of 0.25. This is 2.5 (3.6) times higher than the average CTR on the second (third) position. Therefore, our data supports the prevailing belief in the academic literature that optimizing the rank is crucial for improving a website's traffic. However, there is a considerable standard deviation of 0.19 for position one. Outliers on the first position reach up to 1.0 CTR, while a CTR of 0.0 is also possible. This can also be seen in the second and third positions, with a CTR range between 0 and 0.8.

This range translates to a significant variance reiterating the need to correctly predict CTRs, as even minor deviations in CTR can have significant click and revenue implications for

---

<sup>1</sup>  $\rho$  denotes the Pearson correlation coefficient

businesses and their SEO efforts. However, the strong variance also underlines the importance of other explanatory variables to explain strong deviations. Lastly, the results highlight position as the single most important feature, which inherits substantial implications for the modeling: We infer that the *position* subset is crucial for accurate predictions and that it can thus serve as baseline subset to use in models to analyze the predictive power of other subsets.

## 4.2 Keyword characteristics - the role of a keyword's complexity

Among keyword characteristics, the complexity score reveals a particularly noteworthy effect on CTR. For the complexity score introduced in chapter 3, the simplest possible keyword receives a score of 121, with the score decreasing towards a technically infinite negative value for more complex keywords.

Figure 1 shows the average CTR for keywords on position one depending on their complexity. Additionally, it distinguishes between branded<sup>2</sup> and unbranded keywords. The Figure shows that keyword complexity creates mainly two contrasting effects: For more complex keywords consisting of many and/or longer words<sup>3</sup> (score below 80), the average CTR of branded keywords is significantly lower than for unbranded ones. This changes for values with a complexity score of more than 80, where simple branded queries lead to significantly higher average CTRs. For the most simple keywords, branded keyword results, on average, achieve a CTR of more than 50% higher than unbranded keywords. Additionally, a general tendency of decreasing CTR for simpler keywords can be observed for unbranded keywords.

The results imply that optimizing for position one makes sense for branded keywords that are simple and short, where i.e. the brand is likely to be a very prominent component. However, the results also hint that longer and more complex keywords including brand names can even

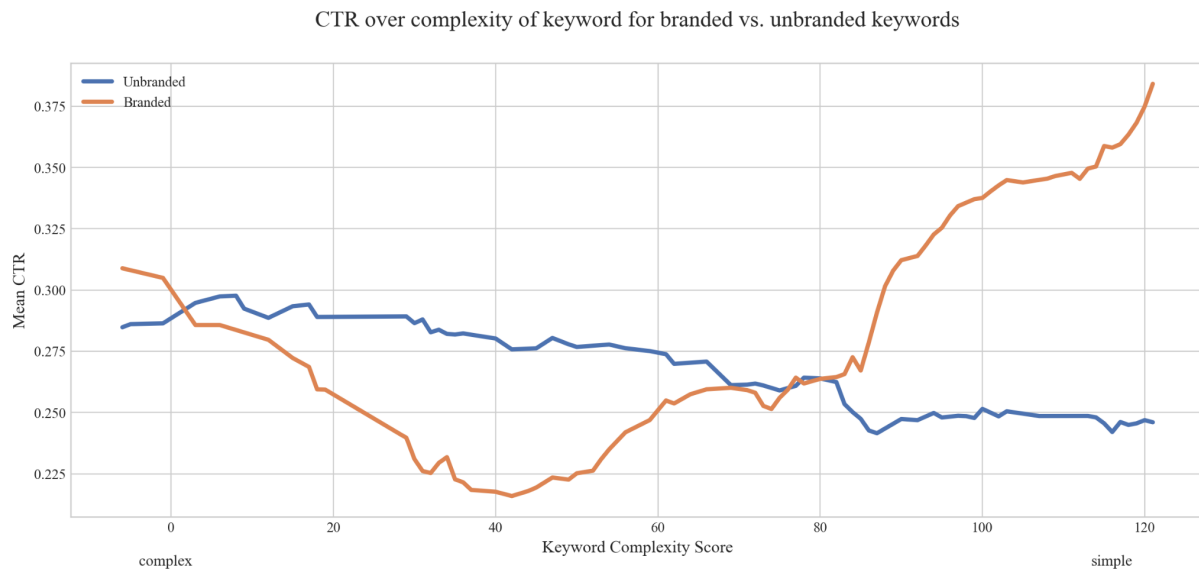
---

<sup>2</sup> A keyword, or search query, is branded if it includes a brand name.

<sup>3</sup> To make the score more tangible: A search query consisting of five words with an average number of syllables per word of 1.8 would achieve a reading ease score of 49.5. An example of such a query is “buy cat food online now”.

harm the CTR of results on position one. This implies that users are more explorative for such keywords and tend to be more likely to consider results after position one.

The results also highlight feature interactions between the complexity of a keyword and whether it is branded, implying a need for prediction models to be capable of modeling feature interactions.



**Figure 1:** Mean CTR over keyword complexity score for branded and unbranded keywords. Mean CTR is calculated as a centered moving average with a total window size of 30. Keywords are more complex for low values and simpler for high values.

### 4.3 SERP features - variation in effects over positions

In the following analysis, we examine the impact of SERP features on CTR and outline implications for SEO decision-making concerning SERP features.

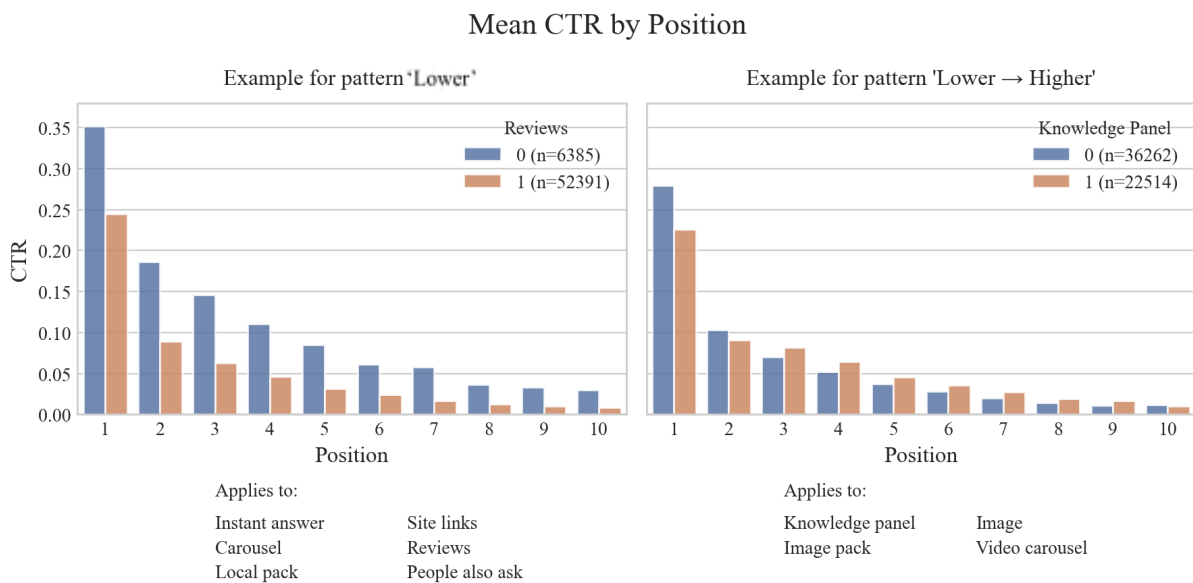
As the count of positional SERP features increases, i.e. a single result has more SERP features bound to it, the CTR tends to rise ( $\rho = 0.22$ ), whereas with more page SERP features, there is a less clear trend, with CTR tending to decrease ( $\rho = -0.11$ ).

In chapter 4.1, we have shown the importance of positions. When adding another dimension, namely the presence of a certain SERP feature, noteworthy tendencies arise. While for some SERP features the average CTR per position remains relatively constant regardless of the



presence of the SERP feature, for others, this presence has a big influence. Tolerating minor exceptions, there are two main patterns observable: With the presence of a SERP feature, first, the average CTR declines over all positions ('Lower'), and second, the average CTR declines for up to the first three positions but increases for subsequent positions ('Lower → Higher'). See Figure 2 for an example of each pattern category and SERP features it applies to.

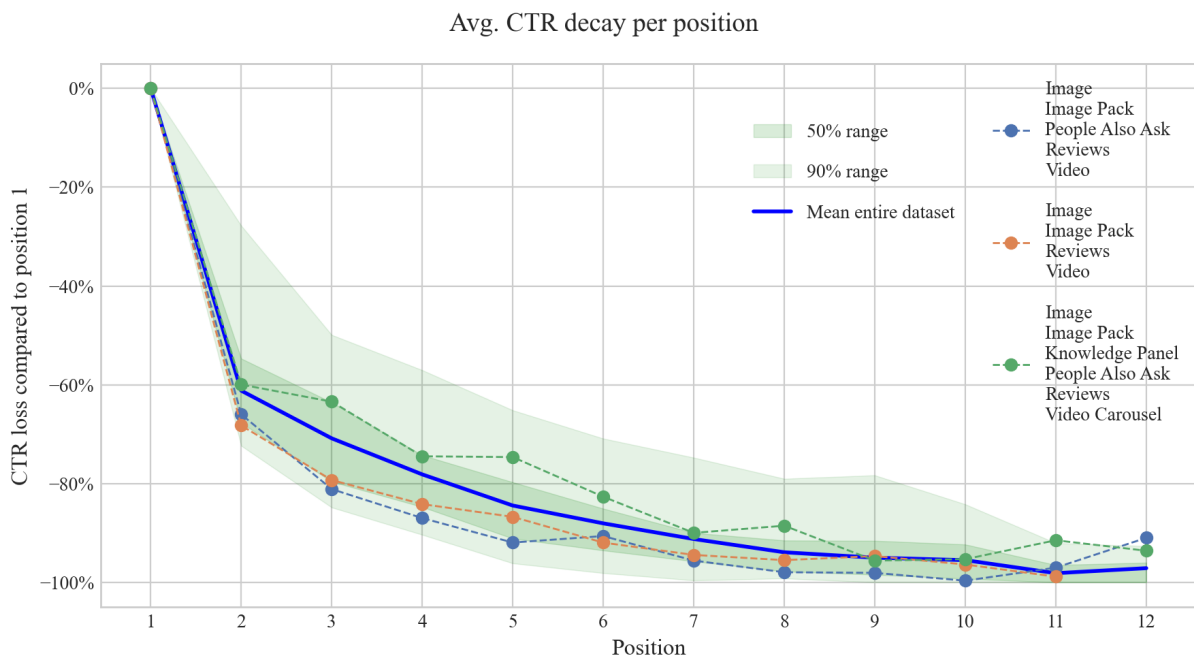
Explanations for why these patterns arise for certain SERP features are not clearly identifiable. However, some features are considerably correlated with intents, e.g. *Image* with transactional search intent ( $p = 0.29$ ). Consequently, when interpreting the results, we must keep in mind that effects might also originate from other feature categories, such as the search intents.



**Figure 2:** Mean CTR per position if a SERP feature is present (1) compared to when it is not present (0). Left plot shows pattern 'Lower' and right plot pattern 'Lower → Higher'. SERP features without a clear pattern or  $n < 1000$  are not listed.

Although we have previously stated that ranking on a top position is of utmost importance for maximizing CTR, the aforementioned findings on SERP features spur the question of whether, in some cases, it is more important to rank on top positions than in other cases. While above, we see the impacts of SERP features isolatedly, in practice, SERP features appear together and thereby likely influence each other. Thus, in the following, we analyze the

importance of positions for each combination of page SERP features. To be able to compare different combinations better, for each combination with more than 200 occurrences, we normalize the CTR to the average value on the first position and then calculate the decay for each subsequent position (Figure 3). As a result, it becomes evident that patterns identified for individual SERP features continue across combinations of them. Among the top three combinations by count, for two, the CTR decays much faster, while for the other one, it decays much slower compared to the entire dataset. For practitioners, this implies that efforts towards ranking on the first position can be of much more value if some SERP features are present, while in the presence of others, a positioning within positions two to five might also suffice. For example, three differing SERP features can already double the maintained CTR on position three (see the difference between the green and blue combination in Figure 3).



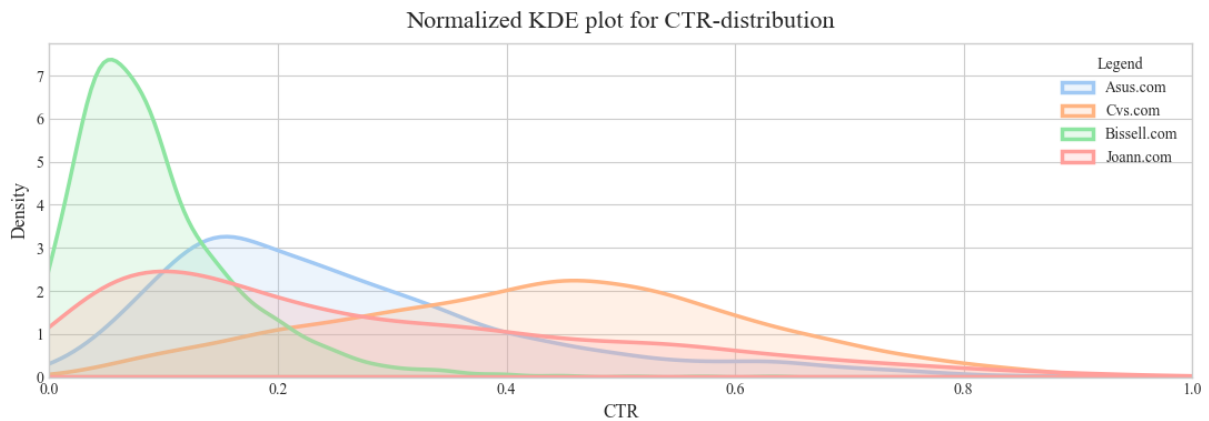
**Figure 3:** Loss in CTR per position. Values are avg. CTR loss for each position in % of avg. CTR on position 1. Only combinations with > 200 occurrences are considered (n=59). Dashed lines represent the top three most frequent SERP feature combinations. Shaded areas represent the range in which 50% and 90% of all values fall.

Together, the implications for the modeling are twofold. First, we show that impactful patterns can be observed on different levels of aggregation, i.e. when using the count of all SERP

features, looking at each feature separately, or looking at combinations of features. Thus, knowledge about SERP features should help models to more accurately predict CTR. Second, we can observe a high degree of interactions. The interactions go beyond the second degree, like the interaction between position and combinations of SERP features. Consequently, to make use of all information inherited in SERP features, models must be able to incorporate feature interactions higher than second degree.

#### 4.4 Result characteristics - how domains & brand perception influence CTR

In this chapter, we focus on the effect of the domains on CTRs. The domains in the analyzed dataset are exclusively from e-commerce websites. They either belong to a manufacturer or a marketplace specialized in product categories like electronics or apparel. Therefore, the domains either represent a brand name or a marketplace and can often be seen as a placeholder for a product category.



**Figure 4:** The Kernel Density Estimation (KDE) plot for four of the top five domains by count, normalized to unify the area under the curve.

Figure 4 shows that domains have widely different CTR distributions on position one. CVS, a well-known pharmacy, achieves the highest average CTR among domains that appear at least 1,000 times. While Bissel, a premium vacuum cleaner company, receives the lowest average CTR of these domains. The kernel density estimated distributions of CTR are right-skewed

for domains with low average CTR, whereas domains with high average CTR go along with a more uniform distribution. We find that there are no correlations with other variables that could clearly explain these patterns. This indicates that the brand name is still important in determining CTR to a website's traffic.

This is especially clear when looking at the top and the worst performing domains on position one, CVS and Bissell. Both of the websites are reached from search queries that include their brand name in 87% of cases for Bissel and 77% for CVS. In both cases, 98% of these queries have a transactional (clear purchasing) intent. However, even though both domains appear on position one, in search queries with a purchasing intent for a product of their brand, there is a large difference in CTR. 'cvs.com' achieves an average CTR of 26%, compared to 6% for 'bissel.com'.

A possible explanation could be a certain brand recognition by the consumers. While the pharmacy, CVS, and its online business is a respected marketplace, the premium vacuum cleaner company, Bissel, is known for its product but not as a primary marketplace, as third parties often offer lower prices. When displayed in Google, Bissel's website is in 70% of the cases enriched with the positional SERP feature *Review*. As the *Review* can also show prices, it is likely that third-party retailers show a lower price and drive traffic away. Additionally, searches for Bissel's products are 55% more likely to include ads than the rest of the data set. This presumably leads to more competition for users' attention.

Supporting the thesis for the high relevance of a website's brand, we find that the two domains with the highest average CTR ('hp.com' & 'cvs.com') are the only ones listed in the top 100 most powerful brands index in the US (Tenet Partners 2020), while the others do not appear.

Based on these findings, we conclude that domains are important in determining the CTR of a website and have high relevance for predicting it.

#### 4.5 Feature interactions - SERP features & the brand name in keywords

In our previous analyses, we have uncovered several feature interactions across and within feature categories. In the following, we quantify the strength of feature interactions and reveal the strongest interaction's underlying forces.

Friedman and Popescu (2008) introduce the H-statistic to measure interactions of two or more features in an ML model. The H-statistic quantifies the interaction strength of two features by the difference in the explained variance of a decomposed partial dependence (PD) function and the observed PD function. The decomposed PD function assumes no interactions, whereas the observed PD function measures interactions. Thereby, it calculates the explained variance resulting from interactions.

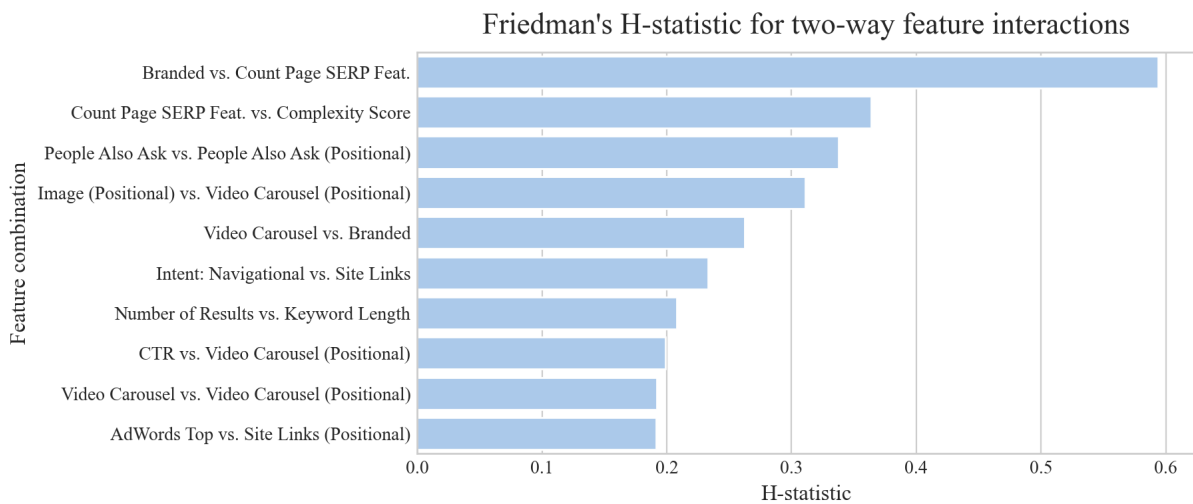
$$(2) \ H_{jk}^2 = \frac{\sum_{i=1}^N [PD_{jk}(x_{ij}, x_{ik}) - PD_j(x_{ij}) - PD_k(x_{ik})]^2}{\sum_{i=1}^N PD_{jk}(x_{ij}, x_{ik})^2}$$

If the resulting score amounts to zero, all variance is explained by the individual 'contributions' of the features. Thus, there is no interaction between them. A higher value indicates that more variance is explained by the PD functions with interactions. This suggests that each individual PD function is constant and the effect on the prediction only comes through the interaction (Friedman and Popescu 2008).

The H-statistic has some limitations that need to be addressed when using it. First, it is non-deterministic, i.e. the same input leads to different results. Therefore, ten iterations were completed to see if it delivers stable results. Second, the H-statistic is very computationally expensive when applied to higher-degree interactions. Thus, we only analyze second-degree interactions. Third, the H-statistic only considers the interaction strength. It does not indicate whether the features have a relevant effect on the target variable or the kind of interaction. For

these reasons, it is necessary to analyze the feature combinations proposed by the H-statistic, to clearly explain their effects. In the following, we point out the most relevant findings.

As Figure 5 reveals, the strongest second-degree interaction of an H-value of around 0.6 appears for the features branded with the count of page SERP features. These features reveal an interesting interaction pattern, which is displayed in Figure 6. Branded keywords perform nearly three times as well as unbranded ones when few SERP features are displayed by Google. This effect reduces on result pages with more SERP features. When more than seven SERP features are present, the effect switches and unbranded keywords perform better than branded ones.

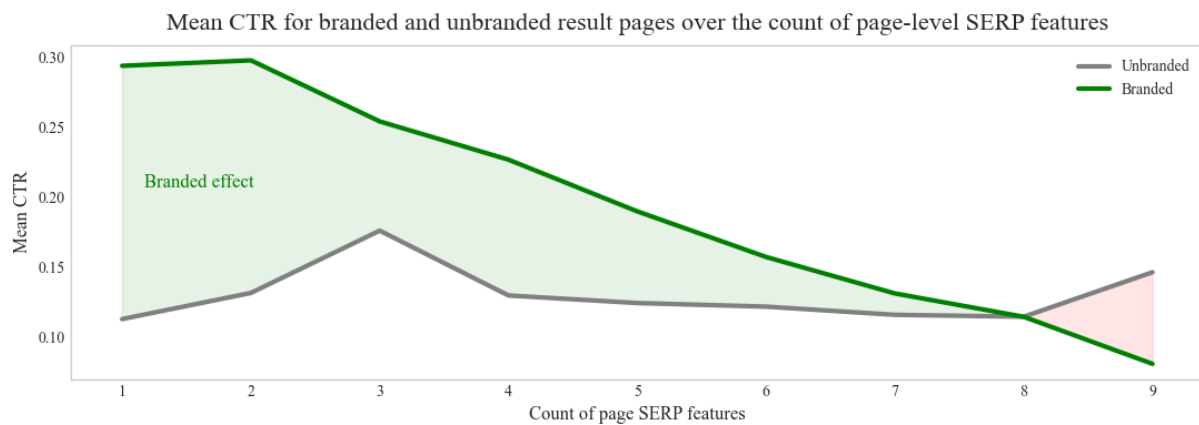


**Figure 5:** Friedman's H-statistic for two-way feature interactions. Highly correlated features and features with less than 500 occurrences were dropped, as they are likely to bias the H-statistic.

As this effect is not intuitively explainable, it is worth investigating it further. We find that result pages with few SERP features have very different properties, depending on whether the keyword is branded or not. The data also shows that result pages for branded keywords are much more likely to have navigational and transactional intent than those for unbranded ones. In contrast, unbranded result pages with few SERP features are more often informational than

their branded counterparts. To summarize, branded result pages with few SERP features have a high CTR, driven by purchase-oriented users who are looking for a specific (branded) website.

The implications are two-fold. First, branded search queries are more likely to generate website traffic if they lead to result pages with few SERP features than to result pages with many SERP features. Second, for unbranded search queries, the number of SERP features is not similarly relevant and a medium CTR can be achieved regardless of them.



**Figure 6:** Mean CTR for branded and unbranded result pages over the count of page-level SERP features.

## 4.6 Findings & implications summary

Next to the detailed analysis of individual drivers of CTR above, there are two main implications based on this data analysis. First, within each subset there are features that have a relevant effect on CTR. Therefore, we will consider all of them as model inputs. Second, we revealed many interactions between features. Most features interact with position, which we either explicitly stated or we analyzed features on only one position, to isolate the feature effects. This shows the high relevance of the position for CTR. Additionally, features within and in between subsets interact with each other. These feature interactions imply that the tested models should either implicitly or explicitly handle feature interactions.

## 5. Modeling

In the following, different models for the prediction of CTR are tested and the predictive power of the four feature subsets is examined. First, we introduce the tested models. Then, the modeling methodology is outlined. Lastly, we analyze the modeling results.

### 5.1 Examined models

The different feature subsets are tested over different models. As pointed out in 2.2, the underlying problem is a regression problem. Consequently, the tested models are regression models. Furthermore, these tested models were selected considering their relevance and performance in recent literature as well as their applicability to the dataset and its characteristics as pointed out in the EDA. Additionally, linear regression models are applied as a benchmark due to their simplicity. Three model categories emerge: Linear regression models, neural network regression models, and tree-based regression models.

#### 5.1.1. Linear regression models

Linear regression attempts to model the relationship between two or more variables by fitting a linear equation to observed data (Uyanık and Güler 2013). This allows for fast computation and easy interpretability. However, linear regression models assume the predictor variables to be independent from one another and linearly related with the independent, predicted variable (Su, Yan, and Tsai 2012). Furthermore, linear regression models are prone to outliers and overfitting the training data in the case of a larger number of features (Filzmoser and Nordhausen 2020). As a result, the model can be simplistic and unable to capture real-life complexities in data. Based on the dataset's characteristics outlined in chapter 3, the limitations also apply to the dataset at hand. There are strong feature interactions, i.e., the predictor variables are not independent, and relations with CTR are not always linear.



Regardless of these drawbacks, we will use ordinary least squares regression (OLS), the most popular linear regression model (Su, Yan, and Tsai 2012), as our baseline model. In an OLS model, the squared errors are minimized (De Gryze, Langhans, and Vandebroek 2007). Using this model allows us to compare more elaborate models with a solid, proven statistical framework. Suppose a model performs worse than OLS linear regression, despite the shortcomings and violated assumptions of linear regression described above. In that case, we can assume that it is not suitable for the dataset at hand. We also adapt the OLS model by adding feature interactions as the product of the values of two features (Poly2). As in this way, the number of features grows immensely, we only allow feature interactions with a pearson correlation coefficient  $> 0.05$  in the training data set to be used by the model. Potential overfitting will be addressed using a third variant, a Ridge regression model, which implicitly performs feature selection to avoid overfitting using L2-Norm.

#### 5.1.2. Tree-based regression models

Decision tree models are based on a series of if-then-else decision rules, resulting in tree-like structures, and are already used since the 1960s (Morgan and Sonquist 1963; Messenger and Mandell 1971; Quinlan 1986). Decision trees are enhanced in state-of-the-art models that use ensembles, i.e. a combination of decision trees to make a joint prediction (Omer and Rokach 2018). The considered tree-based models in this work are all ensemble models based on decision trees. The advantages of such models lie in them not requiring specific distributions in the data, being fast without preprocessing, and they are capable of modeling feature interactions (Agarwal et al. 2022). In addition, on medium and small-sized datasets, tree-based ensemble models often outperform more elaborate state-of-the-art models such as neural networks (Grinsztajn et al., 2022; Hancock and Khoshgoftaar, 2020).

Hence, tree-based ensemble models can address both the interactions and limited sample size

of the available dataset, making them very promising candidates.

We will test the top three performing tree-based models for tabular data, found in a recent comparative benchmark by Grinsztajn et al. (2022), namely Random Forest (Breiman 2001), Gradient Boosting Decision Trees (GBDT) (Friedman 2001) and XGBoost (Chen and Guestrin 2016), as well as CatBoost, another tree-based model achieving high scores in benchmarks (Prokhorenkova et al. 2018).

Random Forests make use of bagging, an approach where the results of multiple independent models, each trained on a subset of the data, are combined (Breiman 1996). Applied to Random Forests, this means that to get a prediction, the results of several randomized decision trees are aggregated through averaging (Breiman 2001). This model is continuously praised for its versatility, speed, and ease to use (Biau and Scornet 2016).

Contrary to bagging models, boosting models build decision trees sequentially and learn from the residuals of their predecessors (Freund and Schapire 1996). A variant of this are *gradient* boosting models, in which each tree, instead of predicting the target itself, predicts the error of its predecessor. Friedman (2001) laid the foundation for this variant with GBDT. Empirical studies indicate that gradient boosting models perform comparatively well on heterogeneous data, but do poorly on homogenous data (Hancock & Khoshgoftaar 2020).

Based on the concept of gradient boosting, Chen and Guestrin (2016) proposed XGBoost, which makes use of advanced regularization techniques (L1 & L2), improving its generalization capabilities. In addition, the training of XGBoost is parallelizable, making it much faster. XGBoost remains the go-to tool for most practitioners and data science competitions (Kossen et al. 2021).

Compared to the other examined gradient boosting models, CatBoost optimizes decision trees for categorical variables through two ways. First, the use of ordered target statistics allows for

efficient handling of categorical features with high cardinality. Second, ordered boosting prevents prediction shift, referring to a phenomenon where what the model learns in the training set is not reflected in the testing set (Prokhorenkova et al. 2018).

### 5.1.3. Neural network regression models

While Deep Neural Networks (DNN) achieve unprecedented performance on tasks related to homogeneous data (e.g. image, audio, and text data), heterogeneous, tabular data is still deemed an “unconquered castle” for DNNs (Borisov et al. 2022; Kadra et al. 2021; Hancock and Khoshgoftaar 2020). However, a recent comparison between the performance of traditional machine learning models and DNNs on tabular data by Borisov et al. (2022), suggests that while for most cases boosting models deliver best performance, on a dataset with more than 10 million samples, DNNs can achieve similar or even better performance. Furthermore, as outlined in chapter 2.2, the best-performing models in recent CTR prediction research incorporate neural networks. Based on these findings, we test a neural network structure especially designed for tabular data, TabNet (Arik & Pfister 2021), as well as Wide&Deep (Cheng et al. 2016) and DeepFM (Guo et al. 2017), neural networks suggested by recent CTR prediction research.

TabNet has been designed for tabular input data and uses the concept of sequential attention to perform row-wise feature selection. This enables it to use its learning capacity only for the most relevant features. Furthermore, TabNet offers a better interpretability than boosting models (Arik and Pfister 2021).

As highlighted in 2.2, CTR prediction research has evolved to currently suggest models that combine a wide and a deep component. These models usually combine the outputs of a deep neural network (deep component) with those of a linear model (wide component) in a common activation function. This model architecture was first proposed by Google as

Wide&Deep (Cheng et al. 2016) and its core strengths lie in its ability to both memorize and generalize (Jais et al. 2019).

State-of-the-art CTR research extends the idea of wide and deep models and adapts them to the specific needs of ad-CTR prediction. Although being outperformed in specific domains, as for example by MaskNet (Wang et al. 2021) on the Criteo dataset, DeepFM (Guo et al. 2017) is still one of the most commonly applied best-performing models in ad-CTR prediction. DeepFM embeds sparse features and interprets the deep component as a neural network while the wide component incorporates a factorization machine. This allows it to capture both two-dimensional feature interactions in the wide component and higher-dimensional interactions in the deep component. As a result, its advantages lie in the flexibility to, in a unified framework, capture low-level feature interactions explicitly and high-level feature interactions implicitly (Guo et al. 2017). However, it also lacks interpretability and imposes high computational complexity (Yang and Zhai 2022).

The implications for organic CTR predictions on dense search engine data are two-fold: On the one hand, DeepFM is capable of capturing feature interactions very well, which is beneficial. On the other hand, it is highly specialized in dealing with sparse feature vectors, which is not necessarily required in the case of result page data and might harm predictions. Furthermore, neural networks often rely on large samples for model training to generate quality predictions (Alwosheel et al. 2018), which is not the case for the available dataset (only medium-sized; ~60k samples).

## **5.2 Methodology**

### **5.2.1 Subset combinations**

To assess each subset's predictive power, different subsets were tested on the models. Testing a model on a subset combination means that all features of the combined subsets are used for

model prediction. Based on both the literature review in 2.1 and our analysis highlighting the position as the single most important feature, the *position* subset serves as a baseline subset. It is then combined with the other subsets.

First, the *position* subset is combined solely with each of the other subsets such that the additional predictive power of each subset can be assessed isolatedly. The *result* subset, however, inherits a structural limitation: Because it mainly consists of one-hot-encoded domains, it is very specialized to the underlying dataset. As such, if a prediction for a new domain was to be made in a production environment, the one-hot encoded domains could not contribute to making a more accurate prediction. Thus, the *result* subset has very limited generalizability, and scores need to be interpreted cautiously.

Based on this limitation, we first check the predictive power of all generalizable subsets combined, i.e. *position* + *keyword* + *SERP features*. To get an understanding of the full potential, we lastly also include the *result* subset to combine all available meaningful features. As a result, 6 subset combinations are tested on each model presented in 5.1.

### 5.2.2 Evaluation metric

In order to compare the effectiveness of different feature subsets and models, we need a common evaluation metric. All models are evaluated using the Root Mean Squared Error (RMSE), defined as

$$(3) \text{ RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}.$$

This metric allows for the penalization of larger errors which is desirable in the present business case as larger errors can lead to a complete miss investment of SEO efforts. As companies are likely to optimize their SEO marketing for the keyword with the highest

expected clicks, it would be a particular financial loss if this estimation is wrong. In addition to that, the RMSE can be used as a quick optimization metric across all model categories, as opposed to, for example, the mean absolute error (MAE).

### 5.2.3 Hyperparameter tuning

To examine each model's full potential, we conducted hyperparameter tuning where reasonable. The hyperparameters for all linear models were kept constant. Tree-based models were tuned based on the hyperparameter sets suggested by Gorishniy et al. (2021) and Grinsztajn et al. (2022). As the tuning algorithm, Bayesian Optimization (Snoek et al. 2012) with 100 iterations is used since it is reported to perform superior compared to random search (Turner et al. 2021). For the comparison of scores of different hyperparameters, 5-fold cross-validation on the training set was used. Please find the full documentation on tested parameter ranges in appendix IV and the resulting parameters in appendix V.

The training of neural networks was limited to 200 epochs with early stopping after 10 epochs without improvement of the RMSE of the validation set. During the training and evaluation of neural networks, we noticed a large variance in the resulting test scores. To compensate for this, we averaged the scores of 10 different experiments to receive the final test score per neural network based model. However, by doing so, an extremely large number of experiments would need to be performed during hyperparameter tuning, resulting in extraordinarily large computation times<sup>4</sup>. Due to this, we could not conduct in-depth hyperparameter tuning for neural networks. Nevertheless, we performed some manual experiments on hyperparameters based on the domain knowledge gained during the EDA. While for most combinations of features subsets, the altered parameters did not lead to significant improvements, scores for the combination of all subsets were drastically improved

---

<sup>4</sup> Example: For setup [iterations = 100; No. optimized subsets = 6; CV = 5; experiments per score = 5; training time = 30 sec], the optimization would take 5.2 days per model

(see appendix VI). Thus, for all neural networks, the default hyperparameters were used for all combinations of feature subsets except the combination of all subsets.

### 5.3 Results & interpretation

In the following, we analyze the model results (Table 1), first concerning how different models perform, and second considering the predictive power each subset incorporates.

#### 5.3.1 Comparison of model performance

Clear patterns arise both between each model and the categories they are grouped in. Apart from summarizing these patterns, we examine whether the assumptions regarding the applicability of models, made in chapter 5.2, hold true.

On average, linear models perform worse than other categories of models, as expected. As described in 5.1.1., we tested the Poly2 model to check the assumption that model performance improves when we allow interactions between features. This assumption holds true as the Poly2 model significantly outperforms the OLS/baseline model. An exception to this is the case when all feature subsets are used, where the error for Poly2 increases drastically. A potential explanation for this is that with more features, the number of interactions grows quadratically, thus leading to so many features that the linear regression model heavily overfits. Ridge only performs marginally better than OLS, indicating that the OLS model works well with the number of features and overfitting is not an issue.

Although mostly outperforming linear models, models based on a neural network architecture cannot keep up with the performance of tree-based models. While scoring very closely to tree-based models for datasets with few features, such as only position, the differences enlarge with the presence of more features and accompanied complexity (see Table 1).

		RMSE of subset combinations					
		Position	Position + Keyword	Position + SERP Feat	Position + Result	Position + Keyword + SERP Feat	Position + Keyword + SERP Feat + Result
Linear Regression	OLS	0.153	0.145	0.145	0.132	0.137	0.125
	Poly2	0.143	0.129	0.131	0.119	0.122	2.554
	Ridge	0.152	0.144	0.144	0.131	0.137	0.123
Tree Based	Random Forest	0.137	0.108	<b>0.120</b>	0.108	0.100	<u>0.088</u>
	GBDT	<u>0.134</u>	0.109	<u>0.122</u>	0.109	0.108	0.091
	XGBoost	<b>0.133</b>	<u>0.104</u>	<b>0.120</b>	<u>0.106</u>	<b>0.098</b>	<b>0.083</b>
	CatBoost	<u>0.134</u>	<b>0.103</b>	<b>0.120</b>	<b>0.105</b>	<u>0.099</u>	<b>0.083</b>
Neural Network	TabNet	0.137	0.121	0.130	0.111	0.118	0.102
	Wide&Deep	0.137	0.123	0.152	0.110	0.128	0.109
	DeepFM	0.136	0.135	0.137	0.113	0.133	0.114

**Table 1:** Benchmark results on different feature subsets. The top results for each feature subset are **bold**. We also underline the second-best results.

Furthermore, TabNet, a model resulting from research focusing on using neural networks for inference on tabular data, performs better than models resulting from CTR research (Wide&Deep and DeepFM). Additionally, the strength of CTR research models on sparse data and weakness on tabular data is reflected in the scores. For sparse subset combinations (*position + result*) their error is only 4% above the error of tree-based models, which increases to 19% for dense subset combinations (*position + keyword*).

It is important to note that even without hyperparameter tuning, tree-based models achieved consistently lower errors than the neural networks (see appendix VI). Based on this, we assume that the considerably worse performance of neural networks is mainly rooted in the model’s architecture instead of the lack of hyperparameter tuning.

Together, this confirms the previously highlighted inapplicability of models resulting from



CTR prediction research to the dataset due to lack of sparsity. Additionally, the neural networks likely suffer from a relatively small number of observations.

Tree-based models perform clearly the best. It is noticeable that the boosting models mostly perform better than the bagging model. While the random forest performs better than GBDT in half of the feature subsets, it is worse than XGBoost and CatBoost in every subset. XGBoost and CatBoost are the two top-performing models for every feature subset. For two subset combinations, they achieve the same RMSE, and for the remaining subset combinations, one outperforms the other just by 0.001 RMSE. These very similar performances spur the question of which model to use under secondary, non-performance-based factors. XGBoost is better documented and more established than CatBoost, hence facilitating the interpretation, tuning, and understanding of the applied model. For all these reasons, we suggest using XGBoost to predict the organic CTR of websites on Google result pages.

The general optimality of tree-based algorithms is in line with current research. The excellent performance of tree-based models indicates that they are able to capture feature interactions on tabular, heterogeneous data well, while also being able to learn from relatively few samples.

Overall, we observe four main patterns arising across the different models. First, models capturing interactions perform better than those that do not. Second, models designed for the use on tabular data make better predictions than models designed for user-item-specific CTR prediction. Third, the performance differences between models amplify with the presence of more features. Fourth, tree-based models clearly outperform all other model categories.

### 5.3.2 Explanatory power of feature subsets

By training the models on combinations of feature subsets, it is possible to evaluate the explanatory power of subsets and derive implications on the overall relevance of subsets for CTR.

When comparing the improvement in RMSE for additional feature subsets, we focus on the results of XGBoost. Thereby, the comparison is not diluted by outliers created through badly performing models. Additionally, we validate these comparisons by running paired t-tests on the results of the, on average, five best models<sup>5</sup>. Thus, we can determine if the difference in RMSE of feature subsets is statistically significant.

We find that RMSE scores improve when adding each of the other subsets to the *position* subset. Adding the *keyword* and *result* subsets notably improves the RMSE, reducing it by 22.6% and 21.1%, respectively. These subsets offer more explanatory power than the *SERP features* subset which only improves the RMSE by 9.8%. When adding *SERP features* to *position + keyword*, the previously best combination, we only see a slight improvement of 4.9%. The combination of all feature subsets results in the lowest RMSE. It reduces the error of the baseline *position* subset by 37.6% and also improves the error of the *position + keyword + SERP features* subset by 16.2%. All of the mentioned differences are statistically significant according to the paired t-test.

Together, this implies that all feature subsets are relevant, as they decrease error when being added and improve the baseline model that solely relies on position. Additionally, the larger improvements in RMSE for the *keyword* and *result* subset than the *SERP feature* subset, indicate that they have a larger effect on CTRs. However, this observation is only an

---

<sup>5</sup> The T-test was conducted with RMSE scores from the following models: CatBoost, XGB, Random Forest, Gradient Boosting and TabNet. The alpha used is 5%.

indication as feature importances are not quantified from a model perspective. Individual part B solves this, by applying interpretation techniques to XGBoost, which analyze feature importances and interactions to understand the relevance of the subsets for CTR.

Increasing the number of combined feature subsets to three and four subsets improves the model performance with each added subset. However, improved results including the *result* subset must be interpreted with caution due to its highlighted lack of generalizability. To increase generalizability individual part C generates a variety of *result* features based on a result's title, that are generally applicable. This helps further understand the predictive power of results.

## 6. Conclusion

By using a dataset that is novel in its comprehensiveness, this work analyzed drivers of CTR and CTR prediction models. It extended existing SEO research by not only considering the position as a driver of CTR, but also aspects grouped into keyword, result, and SERP feature characteristics. Ultimately, the research questions can be answered as follows:

The analysis confirmed the prevailing SEO literature assumption that position is the most important influence on CTR. However, it was found that also keyword, result, and SERP feature characteristics have a significant impact on CTR. Nonetheless, these impacts are highly dependent on the position of a result and subject to feature interactions.

Furthermore, it was shown that state-of-the-art CTR prediction models are inapplicable for predicting average CTR on Google result pages. Instead, the model analysis revealed that tree-based models, specifically CatBoost and XGBoost, are best suited for the underlying prediction task forecasting CTR. The analysis of model results also revealed that keyword and result characteristics tend to have more predictive power than SERP features.

However, several limitations, which restrict the general applicability of the results, need to be acknowledged. The dataset only covered e-commerce stores in the US, thus being limited to a rather niche segment of the entire Google search spectrum. In addition to that, the analyzed dataset contained only 59k rows. This did not only lead to small sample sizes for some feature combinations so that no statistically significant conclusion could be drawn but also restricted the predictive power of neural networks, which can require large amounts of data to work best. Lastly, limited computation resources did not allow for extensive hyperparameter tuning of neural networks. This potentially prevented those networks from performing better in the model comparison.

Future work could tackle these limitations by applying this methodology to similar data from sectors other than e-commerce. Additionally, building on the analysis above, we recommend future research to attempt confirming the findings on a larger dataset which could also provide a solid base for analyses with neural networks. Furthermore, an enhancement of the data with time series information would allow an analysis of consumer behavior and the dynamics of effects over time. Finally, a similar analysis over different nationalities could provide interesting insights into differences in user behavior and support globally operating website providers.

To extend this analysis the individual parts provide an in-depth study on SERP feature effects on a result level and how these can improve the number of clicks to a website. Additionally, model interpretation techniques are applied to assess the model's applicability, reveal patterns in the data, and quantify the subsets' importance. In addition to that, the analysis of results is extended to also incorporate titles of results.

## **B - Improving & understanding CTR predictions through model interpretation**

### **B.1. Rationales for model interpretability**

Since machine learning algorithms gained more relevance for society and businesses, the necessity to explain why they make certain predictions has emerged for many use cases (Du, Liu, and Hu 2019; Molnar 2020). The need for interpretability is based on different reasons that can be grouped into two categories. On the one hand, interpretability can be driven from an application perspective that examines whether a model bases its decision on desirable reasoning. On the other hand, there is a research perspective that focuses on creating more knowledge about underlying structures in the data.

A prominent example from the application perspective is a self-driving car. In this high-risk area, it is crucial to understand how an algorithm recognizes obstacles from pixel values and initiates a breaking maneuver (Du, Liu, and Hu 2019). For example, understanding that an algorithm detects a cyclist based on two round objects, its wheels, can lead developers to spot potential shortcomings. The case in which saddle pockets partially cover a wheel could be introduced to the training data, to ensure a safer behavior of the algorithm (Molnar 2020). Moreover, interpretation techniques are used to detect biases in machine learning algorithms (Molnar 2020), that are implicitly learned from biases in the data itself, for example, models judging based on race, gender, or other unintended factors (Mehrabi et al. 2021).

Additionally, machine learning interpretability is valuable from a research perspective. Due to their predictive capabilities, machine learning algorithms are increasingly used in non-informatic research domains, such as genomics. Researchers apply machine learning to facilitate and speed up the process of DNA sequencing (Yang et al. 2020). Even though

models can predict the sequence with good accuracy, the knowledge and the biological implications remain disguised (Yang et al. 2020). Therefore, interpretation techniques can help translate this knowledge to a humanly understandable level.

This work will apply interpretability techniques to gain insights from both the application and research perspectives. From an application perspective, it will analyze if the model correctly captures previously discovered structures in the data. Additionally, from a research perspective, this work analyzes how and with which strength features impact the predicted CTR.

This will be achieved by identifying suitable model interpretation methods (chapter B.2) and using them to generate knowledge (chapter B.3). The findings contribute to the research community in three ways. First, by quantitatively comparing the importance of SERP features, keyword characteristics and position for organic CTR. Second, the extensive dataset, coupled with a model-driven approach, allows for validating the relationship between organic clicks and paid advertisement competition on a search engine. Third, the effect of appearing within SERP features on CTR is used to improve predictions.

## **B.2 Methodology and theoretical background**

This chapter introduces the feature subsets and the model which are used for interpretation, discusses requirements of interpretability, and introduces interpretation methods that are relevant to this work.

### **B.2.1 Selection of feature subsets and a model type for interpretation**

This work aims at interpreting a model and underlying data that is most feasible regarding business implementations. Therefore, the used feature subsets should offer the best predictive performance without overfitting. These are features from the subsets *position*, *keyword*, and

*SERP features*. Any information based on a specific domain, would not be available in implementations of the model on unseen data. Hence, the *result* subset is not used here. Additionally, the best performing algorithm for this subset combination, the tuned XGBoost, is chosen and interpreted in the following. The model achieves an RMSE of 0.076 on test data, with the specified features. To generalize interpretations of a model it should have high predictive performance to make sure it correctly captures structures in the underlying data (Molnar 2020). This work assumes that the error is small enough to do so. However, it is still a limitation of this work, as the remaining error could bias the findings.

When optimizing prediction scores and using tree-based algorithms, collinearity of features is generally not seen as a problem. Consequently, in the case of XGBoost, it is expected that performance does not suffer from correlated features (Chen et al. 2020). However, collinearity imposes problems when interpreting models (Chen et al. 2020). First, when correlated features describe a similar reality, its importance for a model's predictions is split among the correlated features (Molnar 2020). This is unfavorable, as researchers are interested in the effect of the underlying reality. Second, some importance measures can be biased by creating unrealistic instances or underestimating the importance of a feature by substituting explanatory capabilities from correlated features (Molnar 2020). Therefore, from groups of highly correlated features, only the feature that has the highest Pearson correlation with the target variable is used for training. For example from *position*, *average position per month*, and *previous position* only *average position per month* is applied in the model.

The original dataset used in the previous part of this work included 80,000 observations from US domains. This data set is now enhanced by additional 800,000 observations from the UK with similar features, which gives the findings a higher validity.

### B.2.2 Requirements of interpretability for generating insights in CTR

For this work there are three requirements for the chosen interpretability techniques. First, the techniques need to be able to compare the global<sup>6</sup> relevance of features for the prediction. Thereby, it is possible to further elaborate on the research question which assesses the relevance of SERP features, keyword features, and positional features. Second, the used techniques should detect feature interactions and indicate whether a contribution to the predicted CTR is positive or negative. With these characteristics, it is possible to identify structures in the data and derive implications for SEO decision making. Third, by applying multiple interpretation techniques a more holistic conclusion can be generated and results are fact-checked. This approach is described by Hall: „Several explanatory techniques are usually required to create good explanations for any given complex model. Users should apply a combination [...] of explanatory techniques to a machine learning model and seek consistent results across multiple explanatory techniques.“ (Hall 2019)

These requirements are jointly satisfied by SHAP, permutation importance and gain, which are described from a theoretical perspective in the following chapters and later applied on the CTR prediction.

#### B.2.2 SHAP – measuring a feature’s local contribution towards the prediction

SHAP (SHapley Additive exPlanations) is an often used framework with great properties for comparing the impact of features in a machine learning model (Scott and Lee 2017). It is based on Shapley explanations, which is a game theoretic approach that assigns individual players a share of an aggregated outcome (Winter 2002). This approach can be translated to machine learning by assigning features a contribution to the combined prediction. It allows the following interpretation: “Given the [...] set of feature values (for an observation), the

---

<sup>6</sup> across all observations



contribution of a feature [...] to the difference between the actual prediction and the mean prediction is the estimated Shapley value” (Molnar 2020). To rank features according to their average impact on the prediction, the SHAP importance is calculated by averaging the absolute Shapley values for a feature across all observations.

SHAP is the only additive<sup>7</sup> framework for analyzing feature contributions that obeys the following properties of (a) local accuracy, (b) missingness, and (c) global consistency (Scott and Lee 2017). This means that SHAP explanations (a) sum up to the actual prediction locally, i.e. per observation, (b) do not attribute any importance to missing features and (c) do not decrease for a feature when the actual importance in the underlying model increases (Scott and Lee 2017).

SHAP is applied as the primary framework to generate insights into the model’s predictions and underlying data, because of the previously described properties, its ability to indicate whether a feature’s contribution is positive or negative, and the potential to analyze interactions. Additionally, permutation importance and average gain are used to validate and enhance results from SHAP.

### B.2.3 Permutation Importance and average gain in XGBoost splits

Permutation importance was originally proposed to measure feature importance in random forests (Breiman 2001), but can be applied more generally as a model-agnostic<sup>8</sup> approach. Permutation importance is calculated by randomly shuffling each feature and measuring the decrease in a model’s performance caused by the permutation, in this case, measured as an

---

<sup>7</sup> Additive feature attributions methods create a linear explanation model that consists of binary variables (Scott and Lee 2017).

<sup>8</sup> Model-agnostic methods only examine the input and output space, but do not analyze the internal structure of a model (Du, Liu, and Hu 2019). Therefore, they can be applied to any model, which makes them more flexible and there is no need to restrict a model’s complexity when achieving interpretability (Ribeiro, Singh, and Guestrin 2016).

increase in RMSE. Features for which a large decrease in performance can be registered, are features the model relies on predominantly to make predictions. (Molnar 2020)

Additionally, this work applies the average gain of features in XGBoost trees, which is a model-specific<sup>9</sup> importance measure. It is calculated based on the intrinsic structure of the XGBoost model by averaging the improvement of the loss function in all splits a feature is used in (Friedman, Hastie, and Tibshirani 2013). According to the gain metric, features that have a larger average contribution to reducing the loss function, are more relevant for the overall prediction.

To compare SHAP, permutation and gain importance, they are scaled to have a sum of one across all features. For example, whereas the standard SHAP importance expresses the average absolute contribution of a feature to the prediction, the scaled variant specifies the relative contribution of that feature compared to all other features.

### **B.3 Model interpretation**

#### **B.3.1 In-depth application of SHAP on CTR prediction model**

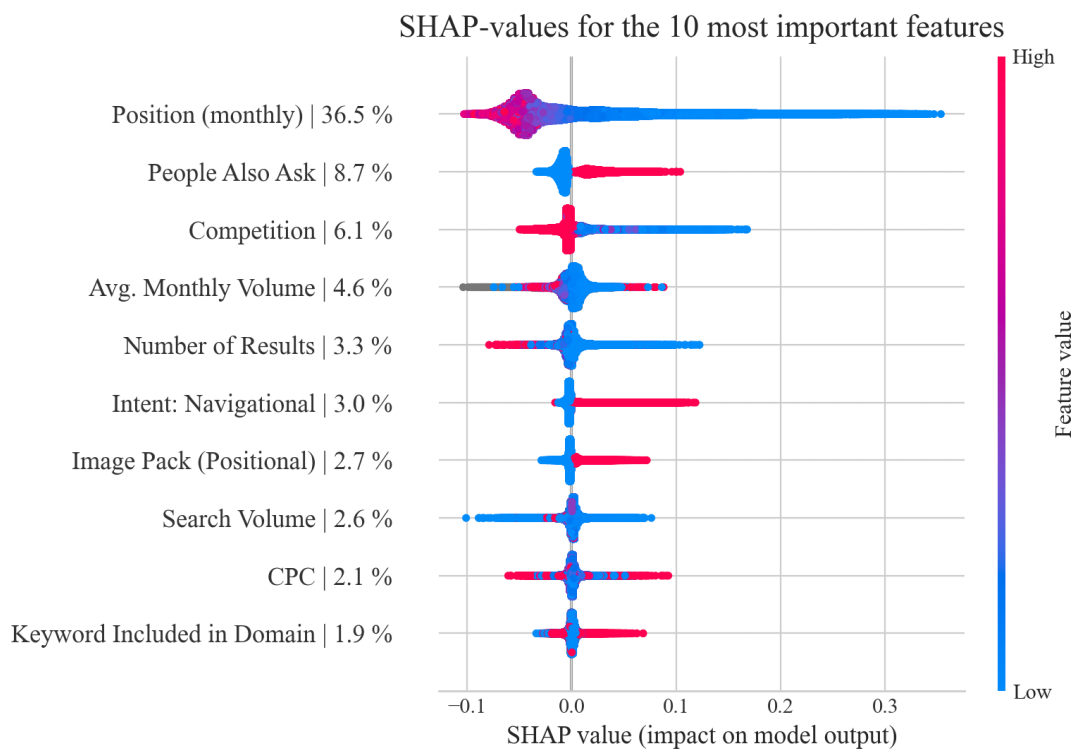
For the following SHAP values were calculated to analyze the XGBoost model's predictions. First, SHAP importances are used to get an overall understanding of the average impact of features on the prediction. Second, certain features and their SHAP values are highlighted to analyze interactions and discover further drivers of CTR.

---

<sup>9</sup> In contrast to model-agnostic methods, model-specific metrics are only applicable to certain model types as they rely on their structure (Molnar 2020).

## B.3.1.1 General overview &amp; SHAP importances

Figure B.1 summarizes the SHAP importances and values<sup>10</sup> for the features with the highest importance scores. The color-coding in Figure B.1 gives an indication of the linearity or non-linearity between feature values and SHAP values. Some features, such as *people also ask* and *image pack (positional)*, imply a linear relationship as high feature values (color-coded in red) are always connected to a positive contribution of the feature to the predicted CTR (SHAP values larger than 0). In contrast, for features such as *average monthly search volume*, high feature values are both associated with a positive and negative contribution to predicted CTR. This is an indication of a non-linear relationship and an interaction with another variable. Both of these effects are analyzed in the following chapters.



**Figure B.1:** The x position of a dot represents the Shapley value, i.e. the contribution of that feature to the prediction, for one observation.

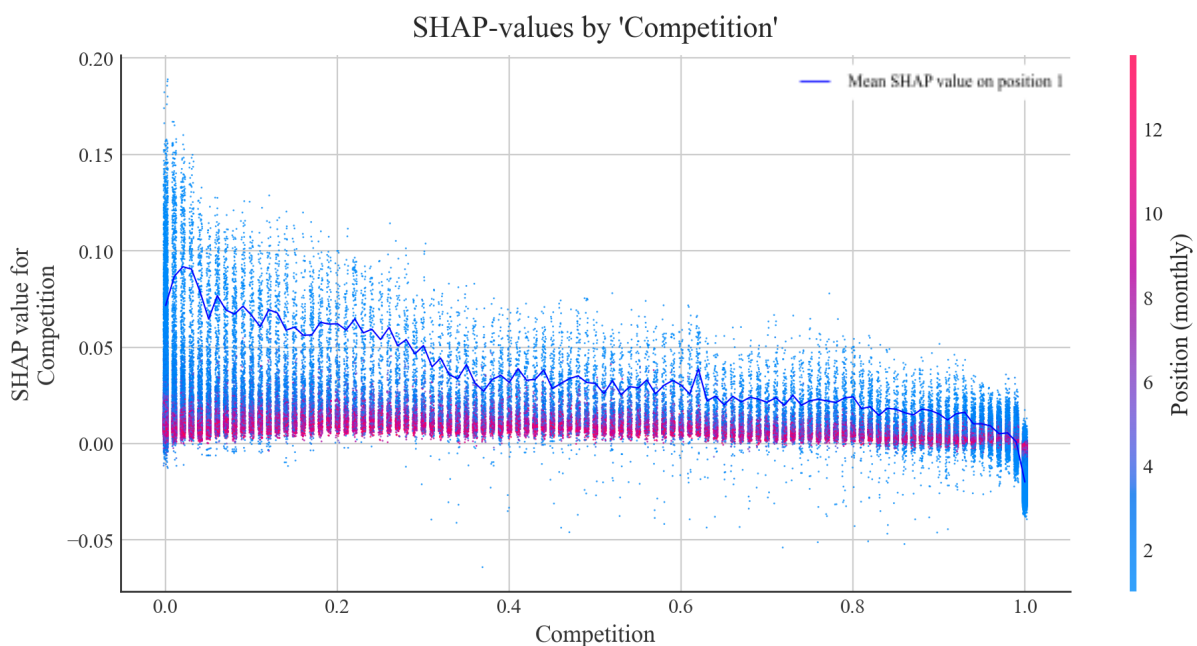
The color indicates the actual feature value. A high density of observations is displayed similar to a violin plot. Next to the feature names, the relative importance of the feature is displayed. For example, averaged across all observations the Position contributes to the prediction with 35.8% of all feature contributions.

<sup>10</sup> As described in B.2.2 SHAP values can be interpreted as “the contribution of a feature [...] to the difference between a specific prediction and the mean prediction” (Molnar 2020).

Additionally, *position* stands out as the most important feature. It has the largest possible contribution to the prediction, as the range of SHAP values is the largest, and has the highest SHAP-importance. Additionally, the SHAP interaction values show interactions of *position* with other features. This is elaborated on in the following for *competition*, but other features like *keyword difficulty*, *keyword count*, *count of SERP features* and *people also ask* show a similar effect.

### B.3.1.2 The relationship between competition in paid advertising and organic clicks

This chapter analyzes the interaction between *competition* and *position* and validates the effect of paid advertisement competition on organic CTR.



**Figure B.2:** Each dot represents one observation. The x-axis shows the feature values for *competition* and the y-axis its Shapley values. The marks are color-coded by average *position*. The blue line highlights the average SHAP value for results on *position* one. As *competition* is not a continuous variable, random noise is added to the x-axis to create visibly identifiable buckets.

A keyword's *competition* only contributes significantly to predictions for websites that are ranked on the top of result pages, which are color-coded in blue in Figure B.2. Contrary, for websites which are on the bottom of a result page (color-coded in red) *competition* only has a

marginal contribution to the prediction. This effect is straightforward, as on high positions the average CTR is higher. Consequently, the absolute contribution of a feature to the prediction, i.e. the SHAP values, should also be higher. Therefore, the model correctly identifies the interaction between *position* and other features, which was previously discussed (see chapters 4.3 & 4.6). This is relevant from an application perspective as the SHAP values reveal that the model behaves as intended.

Moreover, the data in Figure B.2 is valuable from a research perspective, as a low *competition*<sup>11</sup> contributes positively to predicted CTR for results on position one (see the blue line in Figure B.2). This model behavior correctly describes a pattern in the data. For websites on position one the CTR is more than 10 percentage points higher CTR when competition is low, compared to high competition. This effect is less apparent with increasing position values.

While there are several possible explanations for this phenomenon, the following is most logical and supported by findings in the analyzed data set. Keywords with low competition in paid advertising are often associated with an informational search intent (see appendix B.I). In contrast, high competition in paid advertising is more often related to transactional or commercial search intents. Additionally, ads are on average 18 times more frequently displayed for high-competition keywords than for low-competition keywords. As Google Ads are often shown above the first organic results they likely reduce clicks to organic results, based on eye-tracking user behavior (Pernice, Whitenton, and Nielsen 2020). This effect is intuitive and has been reported in Blog articles with unclear data sources (Adlucent 2020; Petrescu 2014). However, to the best of my knowledge, it has never been validated academically and with an extensive dataset, that ads shown on searches with a transactional or commercial intent and high competition in paid advertising reduce clicks to organic results.

---

<sup>11</sup> Competition describes how many advertisers compete for a keyword in paid-search auctions.

Additionally, this finding can potentially be leveraged by businesses in their SEO decision-making process. The high average CTR for low-competition keywords occurs jointly with an above-average search volume. This forms an ideal combination for SEO marketing, because it translates to high website traffic. Making it is especially relevant for companies that leverage organic clicks to create brand and product awareness due to the high informational search intent.

#### B.3.1.3 Improving predictions by explicitly modeling a feature interaction

This chapter describes the interaction between *average monthly search volume* and *people also ask* and points out how understanding the root causes of the interaction improves model performance. The interaction is only noticeable in the UK, which has a share of 88% in the analyzed data. Therefore, this chapter focuses on the UK.

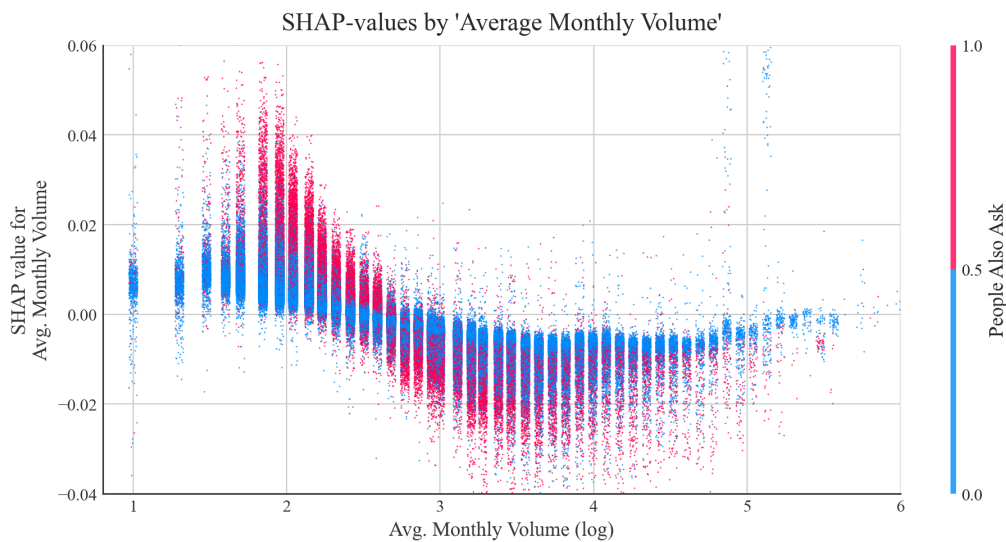
On average, high search volumes have a negative contribution towards the predicted CTR, whereas the opposite is true for low search volumes (see Figure B.3). When *people also ask*<sup>12</sup> is displayed on Google result pages the contribution of *search volume* to the prediction increases. For search volumes of less than 400 searches (below 2.6 on logarithmic scale), the average contribution of *search volumes* to the predicted CTR increases by up to 231% when *people also ask* is present. The negative average contribution of high search volumes is also increased when *people also ask* is present by up to 191%.

At first glance, this interaction is counter-intuitive, as the effect of *people also ask*-boxes on click behavior is not logically linked to search volume. This implies an underlying effect, which can be found in the question of whether a website appears within the people also ask box or next to it. Appendix B.II shows that for low-volume searches, a larger share of the

---

<sup>12</sup> Google's SERP feature *people also ask* shows questions and answers that are similar to the search query. The answers show a text snippet from a third-party website and a link to the site itself.

websites in the analyzed data set appear in the *people also ask*-box. This likely increases CTR as *people also ask*-boxes are a visually prominent element of result pages. In contrast, for high-volume searches, other websites appear more often in the *people also ask*-box, which likely draws clicks away, and the effect of the interaction on CTR is negative (see Figure B.3).



**Figure B.3:** Each dot represents one observation. The x-axis shows the feature values for  $\log(\text{avg. monthly search volume})$  and the y-axis its Shapley values. The marks are color-coded by *people also ask*. As *monthly search volume* is not a continuous variable, random noise is added to the x-axis to create visibly identifiable buckets. Some outliers are cut and data is filtered to the UK.

This has two implications. First, the advantages of appearing in SERP features described in SEO blogs (Moz 2022; Wheelhouse 2022) and by Google (2022b) are validated with a data and model-driven approach. Second, through analyzing XGBoost predictions with SHAP values, a potential for model improvements is revealed. By introducing a feature to the training data that explicitly models whether another website or the website itself appeared in the *people also ask*-box, the RMSE can be improved by 2.3%.

### B.3.2 Comparison of importance metrics and feature subsets

The three previously introduced metrics gain, permutation importance, and SHAP importance were applied to compare the relevance of the feature subsets defined in chapter B.2.

Feature subset / Importance Measure	Gain	Permutation importance	SHAP importance
Position	<b>44.7%</b>	<b>55.9%</b>	<b>36.7%</b>
Keyword	<u>35.9%</u>	<u>27.9%</u>	<u>34.4%</u>
SERP features	19.3%	16.2%	28.9%

**Table B.1:** Subset importance by metric. The most important subset per metric is **bold**, the second-most important subset is underlined.

Table B.1 confirms previous findings, because the subset *position* has the highest importance score across all metrics, followed by the *keyword* and *SERP feature* subset. When attributing feature importances, SHAP importance is most balanced compared to the other metrics, as all subsets have an importance value of 28.9% to 36.7%. Permutation importance splits the importance most diverging, as *position* receives 55.9% and *SERP features* only 16.2%. This difference can be explained by the underlying calculations. As permutation importance uses RMSE it penalizes large deviations over proportionally. Likewise, the average gain is based on the squared loss and thereby has a similar bias towards large outliers. This ‘benefits’ position as it has the largest impact and sets a range of realistic CTR values. For example, permuting the position from 1 to 10 will lead to a different range of expectable CTRs. In contrast, permuting the presence of a single SERP feature has a much smaller effect on RMSE.

However, the interpretation of SHAP values indicates that, on average, SERP features contribute 28.9% to the difference between a prediction and the average prediction. This suggests a larger significance of SERP features for CTR than originally anticipated. These data and model driven findings, should be validated by future work that could analyze behavior of individual users. Additionally, the significance of SERP features could be investigated for non-e-commerce searches.



## B.4 Conclusion and the look ahead

To conclude, from a research perspective this work extends (i) previous findings regarding the importance of feature subsets, and points out (ii) further factors that determine organic CTR. Additionally, from an application perspective it shows (iii) how the model predicts CTR and discovers potentials for model improvement.

First, all importance metrics attribute the same order of importance to feature subsets. The *position* subset is most important for determining CTR, which validates SEO's focus on improving the rank as the central approach for increasing CTR. However, the noteworthy importance attributed to *keyword* and *SERP features* subsets by SHAP importance, indicates that other drivers significantly influence CTR. Therefore, a more holistic view can improve understanding of how users click, and websites can improve traffic. Second, the data and model interpretation approach verifies the relationship between paid advertisement *competition* and organic clicks. This can support practitioners in making better decisions when creating targeted SEO campaigns. Third, the analysis of SHAP values indicates that the XGBoost model recognizes interactions between variables as anticipated, for example, between *position* and *competition*. However, it also revealed how explicitly modeling the appearance of websites in SERP features improves predictions.

Due to the scope of this work only one interpretation method (SHAP) is applied to the fullest and validated by permutation importance and gain. Future research could apply interpretation techniques such as LIME, surrogate models (decision trees), RuleFit, individual conditional expectation (ICE), and partial dependence plots to gain a better understanding of the domain. Additionally, findings could be validated with a user-driven approach, i.e. through eye- or mouse-tracking studies.

## References

- Agarwal, Abhineet, Yan S. Tan, Omer Ronen, Chandan Singh, and Bin Yu. 2022. "Hierarchical Shrinkage: Improving the Accuracy and Interpretability of Tree-Based Methods." *Proceedings of the 39th International Conference on Machine Learning* 162:111-135.
- Alwosheel, Ahmed, Sander van Cranenburgh, and Caspar G. Chorus. 2018. "Is your dataset big enough? Sample size requirements when using artificial neural networks for discrete choice analysis." *Journal of Choice Modelling* 28:167-182.
- Arik, Sercan O., and Tomas Pfister. 2021. "TabNet: Attentive Interpretable Tabular Learning." *Proceedings of the AAAI Conference on Artificial Intelligence* 38, no. 8 (May): 6679-6687.
- Arthur, David, and Sergei Vassilvitskii. 2006. "k-means++: The Advantages of Careful Seeding." *Stanford Technical Report*, (June).
- Bala, Madhu, and Deepak Verma. 2018. "A Critical Review of Digital Marketing." *International Journal of Management, IT & Engineering* 8, no. 10 (10).
- Biau, Gérard, and Erwan Scornet. 2016. "A random forest guided tour." *TEST* 25 (April): 197-227.
- Borisov, Vadim, Tobias Leeman, Kathrin Seßler, Johannes Haug, Martin Pawelczyk, and Gjergji Kasneci. 2022. "Deep Neural Networks and Tabular Data: A Survey." *arXiv:2110.01889v3*, (June).
- Breiman, Leo. 1996. "Bagging predictors." *Machine Learning* 24 (August): 123-140.
- Breiman, Leo. 2001. "Random forests." *Machine Learning* 45, no. 1 (October): 5-32.
- Chen, Hengrui, Hong Chen, Zhizhen Liu, Xiaoke Sun, and Ruiyu Zhou. 2020. "Analysis of factors affecting the severity of automated vehicle crashes using XGBoost model

- combining POI data.” *Journal of advanced transportation*.  
<https://doi.org/10.1155/2020/8881545>.
- Chen, Junxuan, Baigui Sun, Hao Li, Hongtao Lu, and Xian-Sheng Hua. 2016. “Deep CTR Prediction in Display Advertising.” *In Proceedings of the 24th ACM international conference on Multimedia* 16:811–820. <https://doi.org/10.1145/2964284.2964325>.
- Chen, Tianqi, and Carlos Guestrin. 2016. “XGBoost: A Scalable Tree Boosting System.” *KDD '16: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, (August), 785-794.
- Cheng, Heng-Tze, Levent Koc, Jeremiah Hermsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, et al. 2016. “Wide & Deep Learning for Recommender Systems.” *DLRS 2016: Proceedings of the 1st Workshop on Deep Learning for Recommender Systems*, (September), 7-10.
- Chevalier, Stephanie. 2022a. “Global retail e-commerce sales 2026.” Statista. Accessed October 26, 2022.  
<https://www.statista.com/statistics/379046/worldwide-retail-e-commerce-sales/>.
- Chevalier, Stephanie. 2022b. “Global online e-commerce traffic by source and medium.” Statista. Accessed October 28, 2022.  
<https://www.statista.com/statistics/820293/online-traffic-source-and-medium-e-commerce-sessions/>.
- Cui, Meng, and Songyun Hu. 2011. “Search Engine Optimization Research for Website Promotion.” *2011 International Conference of Information Technology, Computer Engineering and Management Sciences*, (09). 10.1109/ICM.2011.308.
- Das, Subhankar. 2021. *Search Engine Optimization and Marketing: A Recipe for Success in Digital Marketing*. N.p.: CRC Press/Taylor and Francis Group.

- De Gryze, Steven, Ivan Langhans, and Martina Vandebroek. 2007. "Using the correct intervals for prediction: A tutorial on tolerance intervals for ordinary least-squares regression." *Chemometrics and Intelligent Laboratory Systems* 87:147-154.
- De Mooij, Marieke. 2017. "Comparing dimensions of national culture for secondary analysis of consumer behavior data of different countries." *International Marketing Review* 34 (3).
- Deng, Wei, Junwei Pan, Tian Zhou, Deguang Kong, Aaron Flores, and Guang Lin. 2021. "DeepLight: Deep Lightweight Feature Interactions for Accelerating CTR Predictions in Ad Serving." *WSDM '21: Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, (March), 922-930.  
<https://doi.org/10.1145/3437963.3441727>.
- Du, Mengnan, Ninghao Liu, and Xia Hu. 2019. "Techniques for Interpretable Machine Learning." *Communications of the ACM* 63 (1): 68-77.  
<https://arxiv.org/pdf/1808.00033.pdf>.
- "The Effect of Paid Search on Organic Traffic." 2020. Adlucent. Accessed November 27, 2022.  
<https://www.adlucent.com/resources/blog/the-effect-of-paid-search-on-organic-traffic/>.
- Filzmoser, Peter, and Klaus Nordhausen. 2020. "Robust linear regression for high-dimensional data: An overview." *WIREs Computational Statistics* 13, no. 4 (July).
- Flesch, Rudolph. 1948. "A new readability yardstick." *Journal of Applied Psychology* 32 (3): 221 - 223.
- Freund, Yoav, and Robert E. Schapire. 1996. "Experiments with a New Boosting Algorithm." *Machine Learning: Proceedings of the Thirteenth International Conference*.

- Friedman, Jerome H. 2001. "Greedy function approximation: A gradient boosting machine." *Annals of Statistics* 29, no. 5 (October): 1189-1232.
- Friedman, Jerome H., Trevor Hastie, and Robert Tibshirani. 2013. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. N.p.: Springer New York.
- Friedman, Jerome H., and Bodgan E. Popescu. 2008. "Predictive Learning via Rule Ensembles." In *Annals of Applied Statistics*, 916-954. 2nd ed. Vol. 3. <http://www.jstor.org/stable/30245114>.
- Gharibshah, Zhabiz, Xingquan Zhu, Arthur Hainline, and Michael Conway. 2020. "Deep Learning for User Interest and Response Prediction in Online Display Advertising." *Data Science and Engineering* 5:12-26.
- Google. 2022a. "Click through rate definition." Clickrate (CTR) - Google Ads-Support. Accessed September 20, 2022. <https://support.google.com/google-ads/answer/2615875?hl=de>.
- Google. 2022b. "Enable Search result features for your site | Documentation." Google Developers. Accessed November 14, 2022. <https://developers.google.com/search/docs/appearance/search-result-features>.
- Google. 2022c. "Data preprocessing for ML: options and recommendations | Cloud Architecture Center." Google Cloud. Accessed November 13, 2022. <https://cloud.google.com/architecture/data-preprocessing-for-ml-with-tf-transform-pt1?hl=de>.
- Google. 2022d. "List of Google Search ranking updates." Accessed November 3, 2022. <https://developers.google.com/search/updates/ranking>.
- Gorishniy, Yury, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. 2021. "Revisiting Deep Learning Models for Tabular Data." *Advances in Neural Information Processing Systems* 34 (*NeurIPS 2021*), (May).

- Graepel, T., J. Q. Candela, T. Borchert, and R. Zerbrich. 2010. "Web-scale bayesian click-through rate prediction for sponsored search advertising in microsoft's bing search engine." *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, 13-20.
- Grinsztajn, Léo, Edouard Oyallon, and Gaël Varoquaux. 2022. "Why do tree-based models still outperform deep learning on tabular data?" *arXiv:2207.08815*, (July).
- Grips. 2022. "Grips - About Us." Grips - Transaction Intelligence for eCommerce. Accessed November 24, 2022. <https://gripsintelligence.com/about>.
- Guo, Huifeng, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. "DeepFM: A Factorization-Machine based Neural Network for CTR Prediction." *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*, 1725–1731.
- Gupta, Sudobh, Neha Agrawal, and Sandeep Gupta. 2016. "A Review on Search Engine Optimization: Basics." *International Journal of Hybrid Information Technology* 9 (5): 381-390.
- Hall, Patrick. 2019. "On the Art and Science of Explainable Machine Learning: Techniques, Recommendations, and Responsibilities." *Proceedings of the KDD 19*.
- Hancock, J.T., and T.M. Khoshgoftaar. 2020. "CatBoost for big data: an interdisciplinary review." *Journal of Big Data* 7 (94): 1-45.  
<https://doi.org/10.1186/s40537-020-00369-8>.
- Iqbal, Muhammad, Muhammad Noman Khalid, Amir Manzoor, Malik Muneeb Abid, and Nazir Ahmed Shaikh. 2022. "Search Engine Optimization (SEO): A Study of important key factors in achieving a better Search Engine Result Page (SERP) Position." In *Sukkur IBA Journal of Emerging Technologies*, 1-15. Vol 6 ed. Vol. 1.  
<http://journal.iba-suk.edu.pk:8089/sibajournals/index.php/sjcms/article/view/924/309>.

- Jais, Imran K., Amelia R. Ismail, and Syed Q. Nisa. 2019. "Adam Optimization Algorithm for Wide and Deep Neural Network." *Knowledge Engineering and Data Science* 2 (1).
- Kadra, Arlind, Marius Lindauer, Frank Hutter, and Josif Grabocka. 2021. "Well-tuned Simple Nets Excel on Tabular Datasets (NeurIPS 2021)." *35th Conference on Neural Information Processing Systems*.
- Kingsnorth, Simon. 2022. *Digital Marketing Strategy: An Integrated Approach to Online Marketing*. N.p.: Kogan Page.
- Kossen, Jannik, Neil Band, Clare Lyle, Aidan N. Gomez, Tom Rainforth, and Yarin Gal. 2021. "SelfAttention Between Datapoints: Going Beyond Individual Input-Output Pairs in Deep Learning." *arXiv:2106.02584*, (June).
- Krrabaj, Samedin, Fesal Baxhaku, and Dukagjin Sadrijaj. 2017. "Investigating search engine optimization techniques for effective ranking: A case study of an educational site." *2017 6th Mediterranean Conference on Embedded Computing (MECO)*, (June).
- Ledford, Jerri L. 2009. *Search Engine Optimization Bible*. N.p.: Wiley.
- Lewandowski, Dirk, Sebastian Sünkler, and Nurce Yagci. 2020. "The influence of search engine optimization on Google's results: A multi-dimensional approach for detecting SEO." *13th ACM Web Science Conference 2021* 21 (06): 12-20.
- Luh, Cheng-Jye, Sheng-An Yang, and Ting-Li D. Huang. 2016. "Estimating Google's search engine ranking function from a search engine optimization perspective." *Online Information Review*, (April).
- Ma, Chao, Yuze Liao, Yuan Wang, and Zhen Xiao. 2016. "F2M: Scalable Field-Aware Factorization Machines." *30th Conference on Neural Information Processing Systems*.
- MacQueen, J. 1967. "Some methods for classification and analysis of multivariate observations." *Berkeley Symposium on Mathematical Statistics and Probability* 5 (1): 281-297.

- Mehrabi, Ninareh, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Glasten. 2021. "A survey on bias and fairness in machine learning." *CM Computing Surveys (CSUR)* 54 (6): 1-35.
- Messenger, Robert, and Lewis Mandell. 1971. "A Modal Search Technique for Predictive Nominal Scale Multivariate Analysis." *Journal of the American Statistical Association* 67, no. 340 (November): 768-772.
- Molnar, Christoph. 2020. *Interpretable Machine Learning: A Guide for Making Black Box Models Explainable*. N.p.: Leanpub.
- Morgan, James N., and John A. Sonquist. 1963. "Problems in the Analysis of Survey Data, and a Proposal." *Journal of the American Statistical Association* 58, no. 302 (June): 415-434.
- Moz. 2022. "What Is A SERP Feature? Common Types And How To Win Them." Moz. Accessed November 12, 2022. <https://moz.com/learn/seo/serp-features>.
- Olson, Eric M., Kai M. Olson, Andrew J. Czaplewski, and Thomas Martin Key. 2021. "Business strategy and the management of digital marketing." *Business Horizons* 64, no. 2 (04): 285-293.
- Omer, Sagi, and Lior Rokach. 2018. "Ensemble learning: A survey." *WIREs Data Mining Knowl Discov* 8 (4).
- O'Neill, Aaron. 2022. "Japan GDP 1987-2027." Statista. Accessed November 27, 2022. <https://www.statista.com/statistics/263578/gross-domestic-product-gdp-of-japan/>.
- Pan, Junwei, Jian Xu, Alfonso Lobos Ruiz, Wenliang Zhao, Shengjun Pan, Yu Sun, and Quan Lu. 2018. "Field-weighted Factorization Machines for Click-Through Rate Prediction in Display Advertising." *In Proceedings of the 2018 World Wide Web Conference* 18.
- Pernice, Kara, Kathryn Whitenton, and Jakob Nielsen. 2020. *How People Read on the Web: The Eyetracking Evidence*. 2nd ed. N.p.: Nielsen Norman Group.



- Petrescu, Philip. 2014. "How Ads Influence Organic Click-Through Rate On Google." Search Engine Land. Accessed November 29, 2022.  
<https://searchengineland.com/different-types-ads-influence-organic-ctr-google-204676>
- Prokhorenkova, Liudmila, Gleb Gusev, Aleksandr Vorobev, Anna V. Dorogush, and Andrey Gulin. 2018. "CatBoost: unbiased boosting with categorical features." *Advances in Neural Information Processing Systems* 31.
- Quinlan, John R. 1986. "Induction of Decision Trees." *Machine Learning* 1 (March): 81-106.
- Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. 2016. "Model-agnostic interpretability of machine learning." *ICML Workshop on Human Interpretability in Machine Learning*. <https://arxiv.org/abs/1606.05386>.
- Richardson, M., E. Dominowska, and R. Ragno. 2007. "Predicting clicks: estimating the click-through rate for new ads." *Proceedings of the 16th international conference on World Wide Web*, 521-530.
- Ruotsalo, Tuuka, Jaako Peltonen, Manuel J. Eugster, Dorota Głowacka, Patrik Floréen, Petri Myllymäki, Giulio Jacucci, and Samuel Kaski. 2018. "Interactive Intent Modeling for Exploratory Search." *ACM Transactions on Information Systems* 36, no. 4 (October): 1-46.
- Sam-Martin, Kristina. 2020. *Google's featured snippets in the Context of Strategic Content Marketing*. N.p.: E-Book - Kristina SamMartin.
- Schultz, Carsten D. 2020. "Informational, transactional, and navigational need of information: relevance of search intention in search engine advertising." *Information Retrieval Journal* 23:117–135.
- Scott, Lundberg M., and Su-In Lee. 2017. "A unified approach to interpreting model predictions." *Advances in neural information processing systems* 30.

- Semrush. 2021. "Basic docs." Semrush Developers. Accessed September 22, 2022.  
<https://developer.semrush.com/api/v3/analytics/basic-docs/#serp-features>.
- Semrush. 2022a. "A Beginner's Guide to Keyword Search Volume." Semrush. Accessed November 3, 2022. <https://www.semrush.com/blog/keyword-search-volume/>.
- Semrush. 2022b. "What Are SERP Features? An In-Depth Guide." SEMrush. Accessed December 1, 2022.  
<https://www.semrush.com/blog/serp-features-guide/#how-to-tell-if-your-competitors%E2%80%99-sites-have-serp-features>.
- Shih, Bih-Yaw, Chen-Yuan Chen, and Zih-Siang Chen. 2013. "Retracted: An Empirical Study of an Internet Marketing Strategy for Search Engine Optimization." *Human Factors and Ergonomics in Manufacturing and Service Industries* 23, no. 6 (11).
- Snoek, Jasper, Hugo Larochelle, and Ryan P. Adams. 2012. "Practical Bayesian Optimization of Machine Learning Algorithms." *Advances in Neural Information Processing Systems* 25 (NIPS 2012).
- Stevenson, Angus. 2010. *Oxford Dictionary of English*. N.p.: OUP Oxford, 2010.
- Su, Xiaogang, Xin Yan, and Chih-Ling Tsai. 2012. "Linear regression." *WIREs Computational Statistics* 4, no. 3 (April): 275-294.
- Tenet Partners. 2020. "2020 Top 100 Most Powerful Brands," The essential brand rises - adapting to core consumer needs. Ranking the brands. Accessed November 28, 2022.  
<https://www.rankingthebrands.com/PDF/Top%20100%20Most%20Powerful%20Brands%202020,%20Tenet.pdf>.
- Tober, Marcus. 2022. "Zero-clicks Study." Semrush. Accessed December 6, 2022.  
<https://www.semrush.com/blog/zero-clicks-study/>.
- Turner, Ryan, David Eriksson, Michael McCourt, Juha Kiili, Eero Laaksonen, Zhen Xu, and Isabelle Guyon. 2021. "Bayesian Optimization is Superior to Random Search for

- Machine Learning Hyperparameter Tuning: Analysis of the Black-Box Optimization Challenge 2020.” *arXiv:2104.10201*, (April).
- Uyanık, Güliden K., and Neşe Güler. 2013. “A study on mutiple linear regression analysis.” *Procedia - Social and Behavioral Sciences* 106 (December): 234-240.
- von der Osten, Barbara. 2022. “SEO Case Studies: Learn From These 7 Success Stories.” Rock Content. Accessed December 2, 2022.  
<https://rockcontent.com/blog/seo-case-studies/>.
- Wang, Fang, Warawut Suphamitmongkol, and Bo Wang. 2013. “Advertisement Click-Through Rate Prediction using Multiple Criteria Linear Programming Regression Model.” *Procedia Computer Science* 17:803-811.
- Wang, Lih-Werm, Michael J. Miller, Michael R. Schmitt, and Frances K. When. 2013. “Assessing readability formula differences with written health information materials: Application, results, and recommendations.” *Research in Social and Administrative Pharmacy* 9 (5): 503-516. <https://doi.org/10.1016/j.sapharm.2012.05.009>.
- Wang, Zhiqiang, Qingyun She, and Junlin Zhang. 2021. “MaskNet: Introducing Feature-Wise Multiplication to CTR Ranking Models by Instance-Guided Mask.” *Proceedings of DLP-KDD 2021*.
- Wheelhouse. 2022. “How SERP Features Can Negatively Impact Your Business (Updated 2022).” Wheelhouse DMG. Accessed December 9, 2022.  
<https://www.wheelhousedmg.com/blog/important-google-serp-features-and-how-to-optimize-for-them/>.
- Wilcoxon, Frank. 1945. “Individual Comparisons by Ranking Methods.” *Biometrics Bulletin* 1, no. 6 (December): 80-83.
- Winter, Eyal. 2002. “The shapley value.” In *Handbook of game theory with economic applications*, 2025-2054. 3rd ed.

- Yan, Ling, Wu-Jun Li, Gui-Rong Xue, and Dingyi Han. 2014. “Coupled Group Lasso for Web-Scale CTR Prediction in Display Advertising.” *Proceedings of the 31st International Conference on Machine Learning* 32 (2): 802-810.
- Yang, Aimin, Wei Zhang, Jiahao Wang, Ke Yang, Yang Han, and Limin Zhang. 2020. “Review on the application of machine learning algorithms in the sequence data mining of DNA.” *Frontiers in Bioengineering and Biotechnology* 8:1032. 10.3389.
- Yang, Yanwu, and Panyu Zhai. 2022. “Click-through rate prediction in online advertising: A literature review.” *Information Processing & Management* 59, no. 2 (03).
- Yuchin, Juan, Yong Zhuang, Wei-Sheng Chin, and Chih-Jen Lin. 2016. “Field-aware Factorization Machines for CTR Prediction.” *Proceedings of the 10th ACM Conference on Recommender Systems* 16 (09): 43-50.
- Yujian, Li, and Liu Bo. 2007. “A Normalized Levenshtein Distance Metric.” *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29, no. 6 (June): 1091-1095.
- Zhang, Weinan, Jiarui Qin, Wei Guo, Ruiming Tang, and Xiuqiang He. 2021. “Deep learning for click-through rate estimation.” *arXiv preprint arXiv:2104.10584*.
- Zhou, Guorui, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. “Deep Interest Network for Click-Through Rate Prediction.” *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* 18 (07): 1059-1068.
- Ziakis, Christos, Maro Vlachopoulou, Theodosios Kyrkoudis, and Makrina Karagkiozidou. 2019. “Important Factors for Improving Google Search Rank.” *Future Internet* 11, no. 32 (01).

## List of Figures

**Figure 1:** Mean CTR over keyword complexity score for branded and unbranded keywords. Mean CTR is calculated as a centered moving average with a total window size of 24. Keywords are more complex for low values and simpler for high values.

**Figure 2:** Mean CTR per position if a SERP feature is present (1) compared to when it is not present (0). Left plot shows pattern ‘Lower’, middle plot pattern ‘Higher’, and right plot pattern ‘Lower → Higher’.

**Figure 3:** Loss in CTR per position. Values are avg. CTR for each position in% of avg. CTR on position 1. Only combinations with > 200 occurrences are considered (n=59). Dashed lines represent the top three most frequent SERP feature combinations. Shaded areas represent the range in which 50% and 90% of all values fall.

**Figure 4:** The Kernel Density Estimation (KDE) plot for the top six domains by occurrence, normalized to unified area under curve and x-axis limited to [0.0; 1.0] CTR-interval.

**Figure 5:** Friedman’s H-statistic for two-way feature interactions. Highly correlated features and features with small number of occurrences were dropped, as they are likely to bias the H-statistic.

**Figure 6:** Mean CTR for branded and unbranded result pages over the count of page-level SERP features.

**Figure B.1:** The x position of a dot represents the Shapley value, i.e. the contribution of that feature to the prediction, for one observation. The color indicates the actual feature value. A high density of observations is displayed similar to a violin plot. Next to the feature names, the relative importance of the feature is displayed. For example, averaged across all observations the Position contributes to the prediction with 35.8% of all feature contributions.

**Figure B.2:** Each dot represents one observation. The x-axis shows the feature values for *competition* and the y-axis its Shapley values. The marks are color-coded by average *position*. The blue line highlights the average SHAP-value for results on *position* one. As *competition* is not a continuous variable, random noise is added to the x-axis to create visibly identifiable buckets.

**Figure B.3:** Each dot represents one observation. The x-axis shows the feature values for *log(avg. monthly search volume)* and the y-axis its Shapley values. The marks are color-coded by *people also ask*. As *monthly*

*search volume* is not a continuous variable, random noise is added to the x-axis to create visibly identifiable buckets. Some outliers are cut and data is filtered to the UK.

## List of Tables

**Table 1:** Impact of the presence of a SERP feature on the average CTR, split between page and positional SERP features. ‘Lower  $\rightarrow$  Higher’ refers to cases where the CTR within the first three positions is lowered and on subsequent positions higher. SERP features that also appear in Figure X, are underlined. SERP features without a clear pattern or  $n < 1000$  are not listed.

**Table 2:** Benchmark results on different feature subsets. The top results for each feature subset are **bold**. We also underline the second-best result

**Table B.1:** Subset importance by metric. The most important subset per metric is **bold**, the second-most important subset is underlined.

# Appendix

## I. SERP Features Illustration



Illustration of SERP features. For a detailed description and examples of each SERP feature we also recommend visiting the Semrush SERP feature guide: <https://www.semrush.com/blog/serp-features-guide>. Image source: Semrush.



## II. Data Dictionary

The texts in the description fields with data source “Semrush” are taken from Semrush (2021).

Variable	Description	Data Source	Example
<i>Ads</i>	Any type of ad on the page. This can include ads on top of the search, at bottom, or shopping ads. Describes whether the SERP feature appears at least once on the page.	Self-engineered	1
<i>AdWords Bottom</i>	A series of ads (up to 4) that appear at the bottom of the first search results page. Describes whether the SERP feature appears at least once on the page.	Semrush	1
<i>AdWords Top</i>	A series of ads (up to 4) that appear at the top of the first search results page. Describes whether the SERP feature appears at least once on the page.	Semrush	1
<i>AMP</i>	AMP pages (i.e., results marked with the word "AMP" and a gray lightning bolt) shown in search results on mobile devices.	Semrush	1
<i>Avg. Monthly Volume</i>	Search volume averaged over months	Keywordplanner	1312.65
<i>Branded</i>	If the keyword includes a brand name.	Grips	1
<i>Carousel</i>	A row of horizontally scrollable images displayed at the top of search results. Describes whether the SERP feature appears at least once on the page.	Semrush	1
<i>Clicks</i>	Number of clicks on targeted search result	Grips	25
<i>Competition</i>	Measure of competition on search term with 0.0 being the lowest and 1.0 being the highest.	Keywordplanner	0.22
<i>Count Page SERP Feat.</i>	Count of positive (=1) page SERP features for search result page.	Self-engineered	2
<i>Count Positional SERP Feat.</i>	Count of positive (=1) positional SERP features for search result.	Self-engineered	5

Variable	Description	Data Source	Example
<i>CPC</i>	Cost per click of ad entry.	Keywordplanner	2.34
<i>CTR</i>	Click-through rate (clicks divided by impressions).	Grips	0.25
<i>Domain</i>	Domain of the URL of each search result.	Semrush	'jomashop.com'
<i>FAQ</i>	A list of questions related to a particular search that shows up for a particular organic search result. When clicked on, each of the questions reveals the answer. Describes whether the SERP feature appears at least once on the page.	Semrush	1
<i>FAQ (Positional)</i>	A list of questions related to a particular search that shows up for a particular organic search result. When clicked on, each of the questions reveals the answer. Describes whether the result is featured in the SERP feature.	Semrush	1
<i>Featured Images</i>	A collection of images is usually displayed at the top of the SERP if Google considers visual results to be more relevant than text results. Only for mobile devices. Describes whether the SERP feature appears at least once on the page.	Semrush	1
<i>Featured Snippet</i>	A short answer to a user's search query with a link to the third-party website it is taken from that appears at the top of all organic search results. Describes whether the SERP feature appears at least once on the page.	Semrush	1
<i>Featured Snippet (Positional)</i>	A short answer to a user's search query with a link to the third-party website it is taken from that appears at the top of all organic search results. Describes whether the result is featured in the SERP feature.	Semrush	1
<i>Featured Video</i>	A video result to a search query that is displayed at the top of all organic search results. Describes whether the SERP feature appears at least once on the page.	Semrush	1
<i>Flights</i>	A block that displays flights related to a search query. Flight results include information on flight dates, duration, the number of transfers and prices. Data is taken from Google Flights. Describes whether the SERP feature appears at least once on the page.	Semrush	1

Variable	Description	Data Source	Example
<i>Hotels Pack</i>	A block that displays hotels related to a search query. Hotel results include information on prices and rating, and allows users to check availability for certain dates. Describes whether the SERP feature appears at least once on the page.	Semrush	1
<i>Image</i>	An image result with a thumbnail displayed along with other organic search results. Describes whether the SERP feature appears at least once on the page.	Semrush	1
<i>Image (Positional)</i>	An image result with a thumbnail displayed along with other organic search results. Describes whether the result is featured in the SERP feature.	Semrush	1
<i>Image Pack</i>	A collection of images related to a search query that is usually displayed between organic search results. Describes whether the SERP feature appears at least once on the page.	Semrush	1
<i>Image Pack (Positional)</i>	A collection of images related to a search query that is usually displayed between organic search results. Describes whether the result is featured in the SERP feature.	Semrush	1
<i>Impressions</i>	Impressions of the result, including estimates for results on the second search page..	Grips	100
<i>Instant Answer</i>	A direct answer to a user's search query that is usually displayed at the top of organic search results in the form of a gray-bordered box. Describes whether the SERP feature appears at least once on the page.	Semrush	1
<i>Intent: Commercial</i>	Trying to learn more before making a purchase decision (e.g. "Subaru vs. Nissan")	Semrush	1
<i>Intent: Informational</i>	Trying to learn more about something (e.g., "What's a good car?")	Semrush	1
<i>Intent: Navigational</i>	Trying to find something (e.g., "Subaru website")	Semrush	1
<i>Intent: Transactional</i>	Trying to complete a specific action (e.g., "buy Subaru Forester")	Semrush	1

Variable	Description	Data Source	Example
<i>Jobs Search</i>	A number of job listings related to a search query that appear at the top of the search results page. Job listings include the job title, the company offering the job, a site where the listing was posted, etc. Describes whether SERP feature appears at least once on the page.	Semrush	1
<i>Keyword</i>	The searched keyword.	Semrush	'jomashop burberry scarf'
<i>Keyword Complexity Score</i>	Flesch reading ease score (Flesch 1948) of the keyword, with higher values indicating easier-to-read keywords.	Self-engineered	112
<i>Keyword Count</i>	Number of words the search query (keyword) consists of.	Grips	3
<i>Keyword Difficulty</i>	Difficulty to rank for given keyword expressed from 0 to 100	Semrush	87
<i>Keyword Included in Domain</i>	Keyword Included in Domain.	Self-engineered	0
<i>Keyword Included in URL</i>	Keyword Included in URL.	Self-engineered	1
<i>Keyword Length</i>	Number of characters in keyword	Grips	23
<i>Knowledge Panel</i>	Panel on right side of results that often includes images, facts, social media links, and other relevant information to the search query. Describes whether SERP feature appears at least once on the page.	Semrush	1
<i>Knowledge Panel (Positional)</i>	Panel on right side of results that often includes images, facts, social media links, and other relevant information to the search query. Describes whether the result is featured in the SERP feature.	Semrush	1
<i>Local Pack</i>	Embedded Google Maps frame on top of search results. Describes whether SERP feature appears at least once on the page.	Semrush	1

Variable	Description	Data Source	Example
<i>Local Pack (Positional)</i>	Embedded Google Maps frame on top of search results. Describes whether the result is featured in the SERP feature.	Semrush	1
<i>Number of Results</i>	Total number of results for search of keyword	Semrush	456789
<i>People Also Ask</i>	A series of questions that may relate to a search query that appears in an expandable grid box labeled "People also ask" between search results. Describes whether SERP feature appears at least once on the page.	Semrush	1
<i>People Also Ask (Positional)</i>	A series of questions that may relate to a search query that appears in an expandable grid box labeled "People also ask" between search results. Describes whether the result is featured in the SERP feature.	Semrush	1
<i>Position</i>	Position of the result among other search results	Semrush	3
<i>Position (monthly)</i>	Average position in the last month.	Grips	2.45
<i>Position Difference</i>	Position subtracted from previous position, i.e. positive values mean a decrease in position	Semrush	-1
<i>Previous Position</i>	Position last month	Semrush	2
<i>Reviews</i>	Organic search results marked with star ratings and including the number of reviews the star rating is based on. Describes whether SERP feature appears at least once on the page.	Semrush	1
<i>Reviews (Positional)</i>	Organic search results marked with star ratings and including the number of reviews the star rating is based on. Describes whether the result is featured in the SERP feature.	Semrush	1
<i>Search Volume</i>	The search volume of a keyword, cleaned from the original (inflated) value.	Grips	75643

Variable	Description	Data Source	Example
<i>SERP</i>			
<i>Features by Keyword</i>	SERP features listed by ID, before one-hot-encoded.	Semrush	1
<i>Shopping Ads</i>	A row of horizontally scrollable paid shopping results that appear at the top of a search results page for a brand or product search query, and include the website's name, pricing, and product image. Describes whether SERP feature appears at least once on the page.	Semrush	1
<i>Site Links</i>	A set of links to other pages of a website that is displayed under the main organic search result and for brand-related search queries. Describes whether SERP feature appears at least once on the page.	Semrush	1
<i>Site Links (Positional)</i>	A set of links to other pages of a website that is displayed under the main organic search result and for brand-related search queries. Describes whether the result is featured in the SERP feature.	Semrush	1
<i>Top Stories</i>	A card-style snippet presenting up to three news-related results relevant to user's search query, which is usually displayed between organic search results. Describes whether SERP feature appears at least once on the page.	Semrush	1
<i>Top Stories (Positional)</i>	A card-style snippet presenting up to three news-related results relevant to user's search query, which is usually displayed between organic search results. Describes whether the result is featured in the SERP feature.	Semrush	1
<i>Trends</i>	How much interest web searchers have shown in a given keyword in the last 12 months.	Semrush	1
<i>Tweet</i>	A card-style snippet displaying the most recent tweets related to a search query. Describes whether SERP feature appears at least once on the page.	Semrush	1
<i>URL</i>	The url of a google search result	Semrush	'https://www.jomas-hop.com/burberry-4031051.html'

Variable	Description	Data Source	Example
<i>Video</i>	Video results with a thumbnail displayed along with other organic search results. Describes whether SERP feature appears at least once on the page.	Semrush	1
<i>Video (Positional)</i>	Video results with a thumbnail displayed along with other organic search results. Describes whether the result is featured in the SERP feature.	Semrush	1
<i>Video Carousel</i>	A row of horizontally scrollable videos displayed among search results. Describes whether SERP feature appears at least once on the page.	Semrush	1
<i>Video Carousel (Positional)</i>	A row of horizontally scrollable videos displayed among search results. Describes whether the result is featured in the SERP feature.	Semrush	1

### III. Feature Subsets

Subset	Features		
Position	Position	Position Difference	Position (monthly)
Keyword	Number of Results	Intent: Transactional	Keyword Length
	Keyword Difficulty	Competition	Keyword Count
	Intent: Commercial	CPC	Branded
	Intent: Informational	Search Volume	Complexity Score
	Intent: Navigational	Avg. Monthly Search Vol.	
SERP Feat.	Instant Answer	Featured Snippet	Site Links (Positional)
	Knowledge Panel	Image	Reviews (Positional)
	Carousel	Jobs Search	Video (Positional)
	Local Pack	Video Carousel	Featured Snippet (Positional)
	Top Stories	People Also Ask	Image (Positional)
	Image Pack	FAQ	Video Carousel (Positional)
	Site Links	Flights	People Also Ask (Positional)
	Reviews	Knowledge Panel (Positional)	FAQ (Positional)
	Tweet	Local Pack (Positional)	Ads
	Video	Top Stories (Positional)	Count Page SERP Feat.
	Featured Video	Image Pack (Positional)	Count Positional SERP Feat.
Result	Keyword in Domain	Keyword in URL	Domain ( <i>One-Hot-Encoded</i> )



## IV. Tested Hyperparameters

Model	Parameter	Values considered for hyperparameter tuning
Random Forest	Max depth	Categorical: [None, 2, 3, 4], p=[0.7, 0.1, 0.1, 0.1]
	Number of estimators	Integer Log Uniform: 10 $\rightarrow$ 3000
	Criterion	Categorical: ['squared_error', 'absolute_error']
	Max features	Categorical: ['sqrt', 'log2', None, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]
	Min samples split	Categorical: [2, 3], p=[0.95, 0.05]
	Min samples leaf	Integer Log Uniform: 1 $\rightarrow$ 50
	Bootstrap	Categorical: [True, False]
	Min impurity decrease	Categorical: [0.0, 0.01, 0.02, 0.05], p=[0.85, 0.05, 0.05, 0.05]
GBDT	Loss	Categorical: ['squared_error', 'absolute_error', 'huber']
	Learning rate	Real Log Uniform: 0.01 $\rightarrow$ 10.0
	Subsample	Real Uniform: 0.5 $\rightarrow$ 1.0
	Number of estimators	Integer Log Uniform: 10 $\rightarrow$ 1000
	Criterion	Categorical: ['friedman_mse', 'squared_error']
	Max depth	Categorical: [None, 2, 3, 4, 5], p=[0.1, 0.1, 0.6, 0.1, 0.1]
	Min samples split	Categorical: [2, 3], p=[0.95, 0.05]
	Min samples leaf	Integer Log Uniform: 1 $\rightarrow$ 50
	Min impurity decrease	Categorical: [0.0, 0.01, 0.02, 0.05], p=[0.85, 0.05, 0.05, 0.05]
	Max leaf nodes	Categorical: [None, 5, 10, 15], p=[0.85, 0.05, 0.05, 0.05]
XGBoost	Max depth	Integer Uniform: 1 $\rightarrow$ 11
	Number of estimators	Integer Uniform: 100 $\rightarrow$ 1000
	Min child weight	Integer Log Uniform: 1 $\rightarrow$ 100
	Subsample	Real Uniform: 0.5 $\rightarrow$ 1.0
	Learning rate	Real Log Uniform: 1e-5 $\rightarrow$ 0.7
	Col sample by level	Real Uniform: 0.5 $\rightarrow$ 1.0
	Col sample by tree	Real Uniform: 0.5 $\rightarrow$ 1.0
	Gamma	Real Log Uniform: 1e-8 $\rightarrow$ 7.0
	Lambda	Real Log Uniform: 1.0 $\rightarrow$ 4.0
	Alpha	Real Log Uniform: 1e-8 $\rightarrow$ 100.0

Model	Parameter	Values considered for hyperparameter tuning
CatBoost	Max depth	Integer Uniform: 3 $\rightarrow$ 10
	Learning rate	Real Log Uniform: 1e-5 $\rightarrow$ 1.0
	Bagging temperature	Real Uniform: 0.0 $\rightarrow$ 1.0
	L2 leaf regression	Real Log Uniform: 1.0 $\rightarrow$ 10.0
	Leaf estimation iterations	Integer Uniform: 1 $\rightarrow$ 10
TabNet	Number decision steps	Categorical: [3, 5, 10]
	Layer size	Categorical: [8, 16, 64]
	Learning rate	Categorical: [0.01, 0.02]
Wide&Deep	Layer 1 size	Categorical: [64, 128, 256]
	Layer 2 size	Categorical: [64, 128, 256]
	Layer 3 size	Categorical: [None, 64, 128, 256]
	Dropout ratio	Categorical: [0.0, 0.01, 0.05]
	Embedding dimension	Categorical: [4, 16, 32]
DeepFM	Layer 1 size	Categorical: [64, 128, 256]
	Layer 2 size	Categorical: [64, 128, 256]
	Layer 3 size	Categorical: [None, 64, 128, 256]
	Use batch normalization	Categorical: [True, False]
	Dropout ratio	Categorical: [0.0, 0.01, 0.05]
	Embedding dimension	Categorical: [4, 16, 32]

## V. Used Hyperparameters

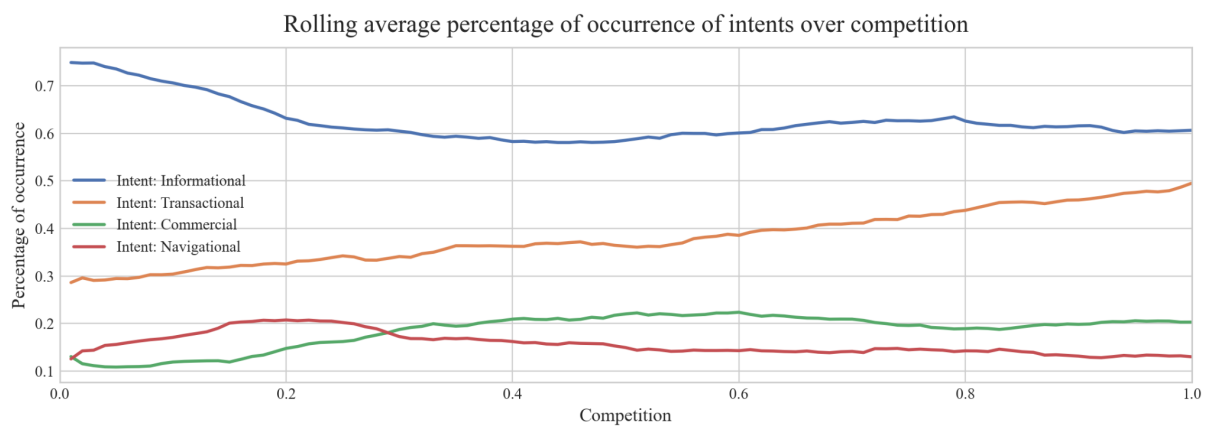
Model	Parameter	Default	Used for Subset					
			Position	Position + Keyword	Position + SERP Feat	Position + Result	Position + Keyword + SERP Feat	Position + Keyword + SERP Feat + Result
Random Forest	Max depth	None	None	None	None	None	None	None
	Number of estimators	100	897	823	3000	3000	175	516
	Criterion	squared_error	squared_error	squared_error	squared_error	squared_error	squared_error	squared_error
	Max features	1.0	0.1	sqrt	0.3	0.4	0.3	0.6
	Min samples split	2	3	2	2	3	2	2
	Min samples leaf	1	38	5	8	9	1	2
	Bootstrap	True	True	False	False	True	False	True
	Min impurity decrease	0.0	0.0	0.0	0.0	0.0	0.0	0.0
GBDT	Loss	squared_error	squared_error	squared_error	squared_error	squared_error	squared_error	huber
	Learning rate	0.10	0.067	0.033	0.030	0.351	0.126	0.767
	Subsample	1.0	0.517	1.0	0.711	0.833	0.634	0.819
	Number of estimators	100	128	497	220	196	487	268
	Criterion	friedman_mse	friedman_mse	friedman_mse	friedman_mse	friedman_mse	friedman_mse	friedman_mse
	Max depth	3	2	None	None	4	5	2
	Min samples split	2	3	2	2	3	3	2
	Min samples leaf	1	24	50	1	10	24	9
	Min impurity decrease	0.0	0.010	0.050	0.050	0.010	0.010	0.010
	Max leaf nodes	None	15	None	15	5	5	15
XGBoost	Max depth	6	11	11	8	11	11	11
	Number of estimators	100	406	689	603	786	460	714
	Min child weight	1	100	1	1	1	1	1
	Subsample	1.0	1.0	0.993	1.0	1.0	1.0	0.668
	Learning rate	0.30	0.015	0.030	0.035	0.016	0.039	0.031
	Col sample by level	1.0	0.50	0.50	0.50	0.50	0.917	0.50
	Col sample by tree	1.0	1.0	1.0	0.50	0.50	0.708	1.0

			Used for Subset					
Model	Parameter	Default	Position	Position + Keyword	Position + SERP Feat	Position + Result	Position + Keyword + SERP Feat	Position + Keyword + SERP Feat + Result
	Gamma	0.0	1.0e-08	7.7e-05	0.009	1.0e-08	1.0e-08	0.001
	Lambda	1.0	4.0	4.0	2.859	1.0	4.0	1.0
	Alpha	0.0	1.0e-08	0.080	0.113	1.0e-08	0.014	1.0e-08
CatBoost	Max depth	6	5	10	9	9	10	9
	Learning rate	0.030	0.013	0.041	0.033	0.026	0.056	0.054
	Bagging temperature	1.0	0.885	1.0	1.0	0.0	1.0	0.0
	L2 leaf regression	3.0	6.601	10.0	10.0	10.0	10.0	10.0
	Leaf estimation iterations	<i>Dynamic</i>	10	10	1	6	10	10
TabNet	Num decision steps	3	3	3	3	3	3	3
	Layer size	8	8	8	8	8	8	64
	Learning rate	0.02	0.02	0.02	0.02	0.02	0.02	0.025
Wide & Deep	Layer 1 size	256	256	256	256	256	256	256
	Layer 2 size	128	128	128	128	128	128	128
	Layer 3 size	64	64	64	64	64	64	64
	Dropout ratio	0.0	0.0	0.0	0.0	0.0	0.0	0.01
	Embedding dimension	4	4	4	4	4	4	16
DeepFM	Layer 1 size	256	256	256	256	256	256	256
	Layer 2 size	128	128	128	128	128	128	128
	Layer 3 size	64	64	64	64	64	64	64
	Use batch normalization	False	False	False	False	False	False	False
	Dropout ratio	0.0	0.0	0.0	0.0	0.0	0.0	0.01
	Embedding dimension	4	4	4	4	4	4	16

## VI. Improvement Through Hyperparameter Tuning

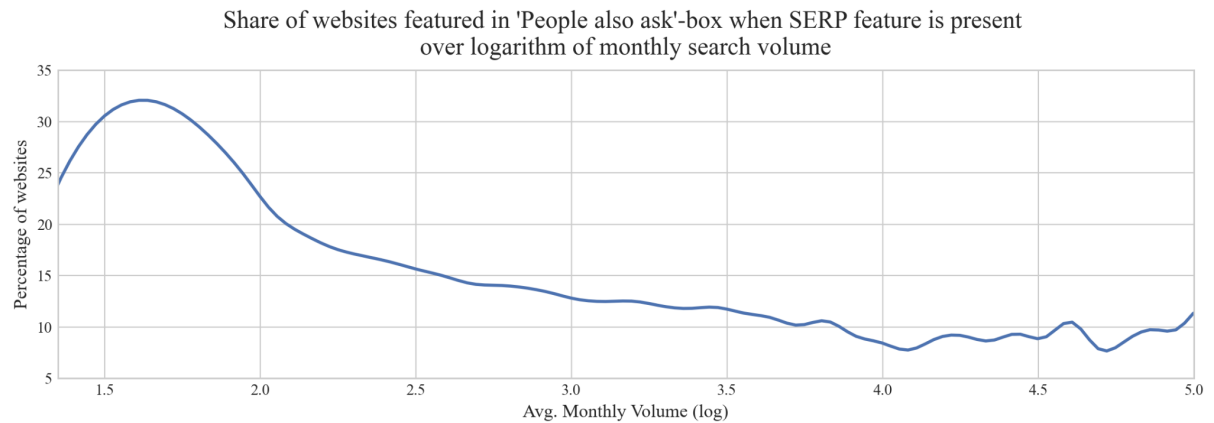
		Improvement over default hyperparameters					
		Position	Position + Keyword	Position + SERP Feat	Position + Result	Position + Keyword + SERP Feat	Position + Keyword + SERP Feat + Result
Linear Regression	OLS	-	-	-	-	-	-
	Poly2	-	-	-	-	-	-
	Ridge	-	-	-	-	-	-
Tree-Based	Random Forest	0.009	0.003	0.007	0.008	0.003	0.01
	GBDT	0.000	0.009	0.002	0.003	0.007	0.007
	XGBoost	0.001	0.005	0.001	0.001	0.004	0.003
	CatBoost	0.000	0.003	0.000	0.000	0.002	0.001
Neural Network	TabNet	-	-	-	-	-	0.033
	Wide&Deep	-	-	-	-	-	0.010
	DeepFM	-	-	-	-	-	0.013

## B.I. Intents by competition



**Appendix B.I:** The graph shows the relative occurrence of search intents by competition. For example, around 60% of keywords with a competition of 0.2 have an informational intent. Note that intents are not ambiguous, i.e. some keywords are mapped to more than one intent. Therefore the percentages do not add up to 100%.

## B.II. Share of websites featured in *People also ask*-box



**Appendix B.II:** The graph shows the percentage of websites that are featured in the *people also ask*-box when this SERP feature is present, i.e. 100% equals to all keywords that have a *people also ask*-box at a given search volume. For example, for keywords with low search volume and *people also ask* present, the website is in the *people also ask*-box in 32% of the cases. In 68 % of the cases another website is in the *people also ask*-box, which likely drives traffic away.