# Leveraging Machine Learning For Optimized Seo In The Healthcare Industry: A Case Study On Everyday Health

## ABSTRACT

Search Engine Optimization (SEO) is critically important in the digital marketing landscape, particularly within the healthcare sector. This thesis develops a robust decision model for search engine rankings aimed at optimizing website performance to better satisfy user demands. Utilizing two advanced gradient boosting models, Light Gradient Boosting Machine (LightGBM) and Extreme Gradient Boosting Decision Trees (XGBoost), this study assesses the relationships and relative importance of various SEO factors. Comparative analysis indicates that XGBoost supersedes LightGBM in predicting actual search engine rankings, achieving an average accuracy rate of 87.7%. A detailed feature analysis by using SHapley Additive exPlanations (SHAP) highlights the significance of internal links, consistent keyword presence across paragraphs, the quantity and length of H2 headings, and the presence of keywords within anchor texts as paramount for effective SEO in the healthcare domain. Conversely, keywords located in footers, URLs, and image alt attributes were found to be less influential. Furthermore, this research includes a practical evaluation of the 'Everyday Health' website through comprehensive webpage crawling. Based on the identified SEO insights, strategic recommendations are provided to enhance the website's search engine positioning. This study not only contributes to the academic understanding of SEO but also offers practical solutions for real-world applications, emphasizing the transformative impact of machine learning in the healthcare sector's digital marketing strategies.

**Keywords:** Search-Engine Optimization; SEO Optimization; Machine Learning; Rank Prediction; Online Healthcare Industry; Digital Marketing; Everyday Health.

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

| No. | Feature | Acronym | Definition |
|---|---|---|---|
| 1 | **Amount of text** | total_words | Count the number of characters in paragraph and titles (<p> and <h> elements) |
| 2 | **H1 count of titles** | h1_num | Count the number of H1 titles on page |
| 3 | **H1 length** | h1_len | Count the average length of H1 titles on page |
| 4 | **H2 count of titles** | h2_num | Count the number of H2 titles on page |
| 5 | **H2 length** | h2_len | Count the average length of H2 titles on page |
| 6 | **H3 count of titles** | h3_num | Count the number of H3 titles on page |
| 7 | **H3 length** | h3_len | Count the average length of H3 titles on page |
| 8 | **Header total** | header_total | Count of all the headers on page |
| 9 | **Image count** | img_count | Count the number of images |

| 10 | **Internal links count** | internalLinks | Count the number of internal links (internal = linking to a page in the same domain) |
|----|--------------------------|----------------|--------------------------------------------------------------------------------------|
| 11 | **External links count** | externalLinks | Count the number of external links (external = linking to a page in the different domain) |
| 12 | **Total links count** | total_link | Count the number of total links (total links = internal links + external links) |
| 13 | **Keyword count H1** | h1_kcount | Count how many times the keyword mentioned in all the H1s |
| 14 | **Keyword count H2** | h2_kcount | Count how many times the keyword mentioned in all the H2s |
| 15 | **Keyword count H3** | h3_kcount | Count how many times the keyword mentioned in all the H3s |
| 16 | **Keyword count p** | p_kcount | Count how many times the keyword mentioned in all the paragraphs |
| 17 | **Keyword in anchor text** | a_kcount | 0 if keyword not in anchor text of any link, 1 if keyword in anchor text of any link |

| 18 | **Keyword in footer** | footer_kcount | 0 if keyword not in footer, 1 if keyword in footer |
|----|----------------------|---------------|----------------------------------------------------|
| 19 | **Keyword in URL** | link_kcount | 0 if keyword not in URL, 1 if keyword in URL |
| 20 | **Keywords in image alt** | imalt_kcount | Count the number of times keyword mentioned in alt tag of images |
| 21 | **Meta desc length** | meta_desc_len | Count the length of the meta description. If no meta description, length = 0 |
| 22 | **Meta keywords count** | meta_kcount | Count the number of meta keywords used |
| 23 | **Page title used** | ti_used | 0 if no page title tag used, 1 if page title tag used |

## I. INTRODUCTION

The rapid advancement of digital technology has led to unprecedented transformations across numerous sectors, with the healthcare industry being one of the most notably impacted. This digital revolution has democratized access to health-related information, thereby empowering patients and consumers to make informed decisions about their health. The widespread availability of internet resources has also intensified the competition among healthcare providers to capture and retain consumer attention online. In this context, Search Engine Optimization (SEO) emerges as an indispensable element of digital marketing strategies aimed at enhancing visibility, user engagement, and ultimately, organizational success.

### The Growing Digitization of Healthcare

The global healthcare market has experienced robust growth, driven by technological advancements and increasing health awareness among populations. According to The Business Research Company, the global healthcare market is expected to grow from $8.45 trillion in 2018 to $11.9 trillion by 2022, reflecting a compound annual growth rate of 8.9% [1]. This growth is mirrored in the digital realm, where health-related searches form a significant portion of internet activity. Google reports that one in every twenty searches is related to health [2], underscoring the critical role of digital platforms in the dissemination of health information.

### The Crucial Role of SEO in Healthcare

Despite the potential of digital platforms, many healthcare providers struggle to navigate the complexities of online marketing. The sheer volume of information and the dynamics of search engine algorithms can be daunting, often leaving valuable content obscured in the vastness of the web. SEO, therefore, becomes a critical tool for healthcare entities to enhance their online presence. Effective SEO strategies can lead to improved search rankings, greater website traffic, and increased patient engagement, which are essential for

healthcare providers in an increasingly competitive market.

A study by Moz in 2021 highlighted that the top five organic search results on Google receive approximately 67.6% of all clicks, with a significant drop in visibility beyond the first page of results [3]. This statistic underscores the importance of a strong SEO strategy; visibility equates to accessibility in the digital age.

**Challenges and Strategic Imperatives**

Despite its importance, the application of SEO in the healthcare industry faces unique challenges. These include maintaining compliance with health information privacy laws, such as the Health Insurance Portability and Accountability Act (HIPAA) in the United States, and ensuring the accuracy and sensitivity of the content, given the potential implications on users' health behaviors and outcomes. Moreover, the rapidly evolving nature of search algorithms necessitates continuous learning and adaptation among digital marketers.

There remains a gap in the specific knowledge and strategic implementation of SEO tailored to the healthcare sector. Many organizations continue to employ generic SEO strategies that may not account for the unique aspects of healthcare services and patient engagement. This thesis aims to bridge this gap by developing a nuanced understanding of SEO factors specifically influential in the healthcare industry and proposing targeted strategies that can lead to more effective web presence management.

**Objective and Structure of the Thesis**

This research is structured to address these complexities through a detailed investigation of SEO effectiveness, employing advanced machine learning models: LightGBM and XGBoost to predict and improve webpage rankings. By analyzing the Everyday Health website, this thesis not only provides theoretical insights but also offers practical, actionable strategies tailored to real-world needs in the healthcare sector.

**The Impact of Internet Resources on Health Decisions**

In 2011, over 80% of adults reported using online resources to assist in making healthcare decisions [4]. As a result, a hospital's website often acts as the initial interaction point for potential users [5]. This makes the hospital or health system's website an essential platform for engaging current and prospective clients, as well as visitors who may accompany patients [6]. The evaluation of a hospital's website by users, and consequently their perception of the hospital, is influenced by their experiences with other consumer websites like Amazon and eBay. Should a hospital's website fail to meet these expectations, it could negatively affect their opinions about the hospital's quality and influence their decision-making [7].

The use of search engines has led to individuals following various paths to their online destinations. Weaver et al. observed that the search behavior of individuals looking for disease-related information differs from those searching for wellness information [8]. In light of these behaviors, many hospital and health system websites have started to feature tools and information that simplify the navigation of complex health situations, enhancing the user experience and promoting a favorable organizational image [9]. Consequently, hospitals are increasingly positioning themselves as reliable advisors, aligning with the accountable care organization (ACO) model, which aims to empower patients and enhance population health outcomes [10].

**Establishing Best Practices in Healthcare Website Design**

The significance of an effective online presence in gaining competitive advantage has prompted researchers to define best practices in website design through standards focusing on accessibility, content, marketing, and technical aspects [11]. Furthermore, the Health Information Technology Institute has set forth specific standards for healthcare websites, encompassing aspects such as credibility, content, disclosure, links, design, interactivity, and caveats [12]. This situation calls for a comparison of US hospital and health system

websites with design standards prevalent in commercial sectors to assess the current landscape.

A high-quality medical or healthcare information website should present comprehensive, unbiased, and accurate information covering both the positives and negatives of the relevant topics. It is also critical that the content is clear and comprehensible to its intended audience. Readability plays a crucial role in the effectiveness of a website, as evidenced by its visitor traffic and its visibility in search engine rankings. Even a website with outstanding content may remain underutilized if it is not accessible or easily understandable to users, getting lost amid the vast array of online information. This study concentrates on websites that are organically ranked by Google, relying on algorithmic selection rather than those boosted through paid advertisements or promotions.

**Addressing Gaps in SEO Research for Healthcare Websites**

While the accuracy and availability of online medical and healthcare information have been extensively explored, considerably less attention has been given to the development of websites dedicated to providing reliable medical information to the public. A search on PubMed, performed on February 22, 2013, using the phrase "development medical information website," yielded only 28 articles that specifically address the creation of medical or healthcare information websites. Search engines like Google, Baidu, and Yandex are frequently used as gateways to access web content, helping users locate the most pertinent information for their needs [13]. These search engines process a vast array of content, crawling through text via links and archiving it in their databases (indexes) for later retrieval and analysis [14]. While the precise algorithms remain proprietary, it is believed that these search engines employ sophisticated computational methods to determine the relevance of web pages to specific search queries [15].

The influence of search engines (SEs) in directing online traffic means that achieving high rankings in search results is highly sought after by website owners. Typically, higher rankings lead to increased site visits and, consequently, higher revenue. This has led

website owners to invest in Search Engine Optimization (SEO), which involves tailoring webpages to align with what are believed to be the ranking factors used by SEs to sort webpages in response to queries. Securing a top position in organic search results is particularly vital for e-commerce businesses, as a substantial share of their traffic often comes directly from search engines [16]. Additionally, e-commerce platforms and other websites frequently utilize blogs to expand and deepen their content, aiming to boost their positions in search rankings. Content quality and relevance are widely recognized as key factors in search ranking by online marketing professionals [17].

**Research Questions and Objectives**

This study aims to address the gap in understanding the impact of content and textual features on search rankings in a practical manner for evaluation by site creators. It seeks to answer the following research questions (RQ):

- RQ1: What is the level of accuracy achieved by gradient models designed to predict webpage ranking in the healthcare industry?
- RQ2: What are the most significant web page ranking factors in the online healthcare industry that can be derived utilizing expert knowledge and machine learning?
- RQ3: What are the top five impactful features and their practical implementations for marketers to enhance web page rankings in the healthcare industry?
- RQ4: How can the Everyday Health website be analyzed to provide marketing and SEO strategic recommendations for upgrading their webpage ranking to the top?

By addressing these research questions, we hope to provide actionable insights for online healthcare industry professionals looking to improve their SEO strategies and increase their visibility and success. We believe that our research has the potential to significantly impact the digital healthcare industry by using machine learning to explore new features that can improve webpage ranking in the education industry through SEO.

## II.   LITERATURE REVIEW

### 1.   Collection of Articles

The impact of ranking factors on web page rankings in search engines has been the subject of extensive research. Search engines typically provide a limited list of ranking factors without disclosing their overall significance in the ranking algorithm. Consequently, numerous studies have aimed to identify key ranking factors and investigate their impact on web page rankings. In this study, we conducted a comprehensive search for empirical research on search engine optimization (SEO) utilizing keywords such as "search engine optimization," "webpages ranking in healthcare industry," "Google rankings," and "using machine learning".

### 2.   Findings forms Related Works

The proliferation of digital technology has significantly impacted various sectors, including healthcare. In the realm of digital marketing, Search Engine Optimization (SEO) has become an essential strategy to enhance the online visibility and engagement of healthcare providers. This literature review provides an in-depth examination of the existing research on the implementation of SEO in healthcare, the challenges encountered, and effective strategies for improving search engine rankings.

SEO is a critical component of digital marketing strategies, designed to improve a website's visibility on Search Engine Results Pages (SERPs). Ledford defines SEO as optimizing various aspects of a website, including content, structure, and technical elements, to enhance its relevance and authority in the eyes of search engines [18]. In healthcare, effective SEO practices ensure that health-related information is easily accessible to users seeking medical advice online. Taneja and Toombs emphasize that high search engine rankings are associated with increased website traffic and greater patient engagement [19]. Chaffey and Ellis-Chadwick further assert that websites on the first page of search results receive the majority of clicks, providing a competitive advantage to those with strong SEO performance [20].

**Strategies for Effective SEO in Healthcare**

Effective SEO strategies in healthcare involve a combination of technical optimization, high-quality content creation, and user experience enhancements. Technical optimization includes improving website speed, mobile responsiveness, and site architecture to ensure that search engines can efficiently crawl and index the site. Patel et al. emphasize the significance of these technical factors in achieving high search engine rankings [21]. The importance of an optimized site structure is echoed by Lee et al., who found that simplified URL structures, internal redirects, and XML sitemaps significantly contribute to better search engine visibility [22].

High-quality, informative, and relevant content is critical for attracting users and signaling search engines about the website's value. Content that addresses common health concerns and provides evidence-based information tends to rank highly on SERPs, as found by Li and Bernoff [23]. Moreover, strategically incorporating keywords within the content can further enhance SEO performance. Zhang and Dimitroff's research highlights that using keywords in both the title and body text results in better performance than using them in just one of the two. Additionally, the study found that using duplicate keywords in the title increases ranking up to a certain point, after which there is a decrease in visibility [24].

7

User experience (UX) is another crucial aspect of SEO. A positive UX, characterized by easy navigation, clear information architecture, and engaging design, reduces bounce rates and increases the time users spend on the site. Websites with superior UX are favored by both users and search engines, leading to better search rankings, as highlighted by Krug [25]. Ensuring a seamless user experience involves removing expired links and content, using canonical URLs, and maintaining a hierarchy within the website structure [22].

**The Role of Metadata and Content Characteristics**

Metadata plays a pivotal role in SEO by providing search engines with information about a webpage's content. Zhang and Dimitroff investigated the impact of metadata on webpage visibility and found that pages with metadata had higher visibility than those without, particularly when the metadata was also included in the page's text content [24]. Their study also revealed that incorporating keywords in both the title and body text enhances search rankings more effectively than using them in only one location [24].

Evans' analysis of Google rankings identified several factors influential for higher rankings, including PageRank score, the number of inbound links, the age of the domain name, and listings in reputable directories [26]. Similarly, Malaga's experimental study emphasized the importance of link building, finding that links from reputable websites significantly impact search rankings [27]. These findings are supported by Wang et al., who identified link popularity as the most critical criterion for high search engine rankings. Their study recommended limiting website title length, optimizing page size, maintaining a hierarchical directory structure, and ensuring appropriate keyword density [28].

**Machine Learning in SEO**

Advancements in machine learning have introduced new opportunities for optimizing SEO in healthcare. Machine learning models, such as LightGBM and XGBoost, can analyze vast amounts of data to identify patterns and predict webpage rankings. Chen and Guestrin demonstrated that these models could enhance the accuracy of SEO predictions, helping

healthcare organizations implement more effective optimization strategies [29]. These models utilize decision trees and boosting techniques to handle large datasets and make accurate predictions based on numerous variables.

Liu et al. applied machine learning to analyze SEO factors, revealing that factors such as page load time, content relevance, and backlink quality significantly impact search engine rankings [30]. Their research demonstrated that machine learning models could effectively identify the most influential factors affecting webpage rankings, providing valuable insights for healthcare providers aiming to optimize their websites.

Zhang and Cabage conducted a comparative study on the effects of link building and social sharing on search rankings, finding that link building had the most substantial impact over 18 months [31]. While social sharing resulted in a rapid increase in traffic, its effect on search rankings was temporary. This study underscores the importance of sustained SEO efforts, particularly in building high-quality backlinks, to achieve long-term improvements in search engine visibility.

The literature highlights the critical role of metadata, content characteristics, and machine learning in optimizing SEO for healthcare websites. Effective strategies combine technical optimization, high-quality content creation, and user experience enhancements to improve search engine rankings.

## III. METHODOLOGY

### 1. Research Context

This research focuses on applying machine learning algorithms to enhance Search Engine Optimization (SEO) in the healthcare industry. The healthcare sector is increasingly dependent on digital platforms to provide information, engage with patients, and offer services. According to a report by Google, one in every twenty searches is related to health, illustrating the critical role of SEO in the healthcare industry [2]. Furthermore, a study by BrightEdge found that organic search drives 53% of all website traffic, making it the most

significant source of web traffic in the healthcare sector [32]. These statistics underscore the growing importance of SEO as healthcare organizations strive to improve their online presence and connect with more patients.

Traditional SEO methods, while effective, can be time-consuming and require substantial expertise, presenting challenges for healthcare providers seeking to maintain a competitive edge. The rapidly evolving nature of search engine algorithms adds to these challenges, necessitating continuous updates and adaptations in SEO strategies. By leveraging machine learning models, this study aims to automate and optimize the SEO process, analyzing large quantities of data to predict webpage ranking patterns. This approach not only streamlines SEO efforts but also provides actionable insights that can significantly enhance a healthcare provider's online visibility.

The healthcare industry is unique due to its stringent regulatory environment, including compliance with laws like the Health Insurance Portability and Accountability Act (HIPAA). These regulations ensure the protection of patient information, which adds complexity to the implementation of SEO strategies. This study addresses these challenges by developing machine learning models that can identify significant features impacting SEO, such as keyword selection, content quality, and backlinking strategies, while maintaining compliance with healthcare regulations.

By focusing on the Everyday Health website as a case study, this research aims to provide practical, actionable strategies for healthcare providers. The models developed in this study - LightGBM and XGBoost are designed to predict webpage rankings with high accuracy, offering insights into the most effective SEO practices for the healthcare sector. The successful implementation of these strategies can help healthcare organizations enhance their online presence, attract more patients, and improve patient engagement.

The healthcare industry is highly competitive and constantly evolving, with a vast array of online resources available to patients and healthcare providers. Effective SEO is crucial for healthcare organizations to maintain their online presence and stay competitive. The

application of machine learning algorithms has the potential to automate and streamline the SEO process, enabling healthcare providers to achieve their SEO goals more efficiently and effectively. Ultimately, the findings of this research can contribute to the development of more effective and efficient methods for optimizing SEO in the healthcare industry, ensuring that valuable health-related information reaches the intended audience.

## 2. Data Collection

To collect data for this research, a comprehensive approach was taken to ensure the relevance and quality of the information gathered. The process began with identifying 150 keywords related to various aspects of the healthcare industry. These keywords were carefully selected to cover a broad range of health topics, ensuring that the data collected would provide a holistic view of the industry's online content. The selected keywords included queries such as "How to maintain a healthy weight," "Best exercises for heart health," "Tips for lowering cholesterol naturally," and "How to boost the immune system," among others. This diverse set of keywords was intended to capture a wide spectrum of healthcare-related searches, from general wellness advice to specific medical conditions and treatments.

To streamline the data collection and analysis, the keywords were clustered into several main topics in following table:

| Cluster | Keywords |
|---------|----------|
| **General Health and Wellness** | How to maintain a healthy weight, Best exercises for heart health, Tips for lowering cholesterol naturally, How to boost the immune system, Natural remedies for better sleep, How to detox the body naturally, Tips for healthy meal planning, How to stay healthy during flu season, Best daily supplements for health, How to improve gut health naturally, Tips for maintaining strong bones, How to prevent common colds |

| | |
|---|---|
| **Chronic Diseases and Conditions** | What are the early signs of diabetes?, How to manage high blood pressure, Symptoms of thyroid problems, Best diet for managing arthritis, How to treat acid reflux at home, What are the best treatments for eczema?, How to relieve migraine pain, Home remedies for urinary tract infections, Early signs of Parkinson's disease, Diet tips for managing Crohn's disease |
| **Skin and Hair Care** | How to have clear skin naturally, Best anti-aging skin care routines, How to reduce acne breakouts, Home remedies for dry skin, Tips for healthy hair and nails |
| **Mental Health** | How to reduce stress quickly, Techniques for managing anxiety, Best therapies for depression, How to deal with panic attacks, Signs of mental health issues, How to overcome social anxiety, Ways to boost your mood naturally, Techniques for improving concentration, Strategies for dealing with grief, How to build emotional resilience, How to handle anxiety without medication, Coping strategies for PTSD, Benefits of cognitive therapy, How to recognize signs of addiction, Managing stress in a high-pressure job |
| **Women's Health** | How to ease menstrual cramps, Pregnancy health tips, How to detect signs of breast cancer, Health tips for menopausal women, Best exercises during pregnancy, How to prepare for a mammogram, Best foods for fertility, How to perform a breast self-exam, Tips for a healthy pregnancy, Managing symptoms of polycystic ovary syndrome (PCOS) |
| **Men's Health** | How to prevent male hair loss, Tips for prostate health, How to increase testosterone naturally, Health screenings every man needs, How to build muscle after 40, How to check for testicular cancer, Ways to prevent balding, How to handle male depression, Health tips for men over 60, Exercises for a healthy prostate |
| **Children's Health** | Vaccination schedule for children, How to deal with common childhood illnesses, Teenage mental health tips, Nutrition advice for growing kids, How to promote healthy eating habits in teens, How to baby-proof your home, Nutritional needs for toddlers, How to handle toddler tantrums, Tips for introducing solid foods, |

| | |
|---|---|
| | How to choose a pediatrician |
| **Senior Health** | How to stay active in old age, Health concerns for seniors, How to improve memory for elderly, Safety tips for seniors living alone, Managing chronic conditions in old age, How to choose a retirement home, Activities for seniors with limited mobility, Best diets for the elderly, How to manage arthritis pain in old age, Safety tips for elderly drivers |
| **Nutrition and Diet** | How to start a ketogenic diet, Benefits of a plant-based diet, Best foods for brain health, How to calculate daily caloric needs, Benefits of omega-3 fatty acids, Vegetarian sources of protein, How to start an anti-inflammatory diet, Gluten-free diet benefits, Best foods to eat in spring for health, How to reduce stress quickly |
| **Fitness and Exercise** | Best morning exercises, How to start running, Yoga poses for beginners, How to lose belly fat, Exercises to improve flexibility, How to train for a marathon, Benefits of Pilates for beginners, How to increase flexibility after 50, Best low-impact exercises for weight loss, How to choose the right fitness tracker |
| **Alternative Medicine** | Benefits of acupuncture for pain relief, What is Reiki and how does it work?, Herbal supplements for anxiety, Benefits of chiropractic adjustments, How to use essential oils safely |
| **Seasonal Health** | How to avoid heatstroke in summer, Tips for winter skin care, How to deal with seasonal allergies, Preparing your immune system for winter |
| **Travel Health** | How to prevent jet lag effectively, Vaccinations needed for South America, Health tips for backpacking in Asia, How to avoid food poisoning while traveling, Essentials for a travel health kit |

| | |
|---|---|
| **General Medical Information** | How to check for skin cancer at home, Importance of annual physical exams, How to lower risk of Alzheimer's, Screening tests every adult should have, Ways to reduce risk of osteoporosis, How to start a meditation practice, Benefits of quitting smoking today, How to reduce screen time effectively, Tips for a balanced work-life, How to improve sleep quality, What causes chronic fatigue?, How to treat acid at home, Remedies for severe headache, Managing symptoms of menopause, Natural treatments for swollen ankles |
| **Pet Health** | How to care for an aging dog, Best diet for cats with kidney disease, Common health issues in parrots, How to prevent ticks on pets, Emergency care tips for pet owners |
| **Technological Health Innovations** | Latest advancements in wearable technology, Benefits of telehealth consultations, How fitness trackers can improve health, Using apps to monitor mental health |
| **Supplements and Vitamins** | Best vitamins for eye health, How to choose the right multivitamin, Benefits of magnesium supplements, Natural sources of vitamin D, Risks of overusing dietary supplements |
| **Special Populations** | Health guidelines for transgender individuals, Managing health with disability, Nutrition tips for the elderly, Fitness routines for wheelchair users, Health considerations for pregnant teenagers |

*Table 1: Keywords Utilized for Crawling Top 10 Webpage Rankings Using Google Search Engine*

Following the determination of the relevant keywords, the study employed the Python library BeautifulSoup to conduct data scraping from webpages. BeautifulSoup was chosen for its simplicity, flexibility, and compatibility with other Python libraries. It offers features that make it easy to extract data from webpages, navigate document structures, and handle errors. Additionally, it can work with various markup languages, including HTML, XML, and XHTML, and integrate with other libraries for HTTP requests and data manipulation [33]. These advantages make BeautifulSoup a preferred tool for scraping datasets from webpages.

To minimize the impact of personalization on search results, the data scraping was carried out in the privacy mode of the Google Chrome browser. This approach helped ensure that the search results were as unbiased as possible. Although search results can vary based on temporal and geographic factors, this study utilized a cross-sectional sample to represent a snapshot of the healthcare industry at a specific moment in time.

For each keyword, the top ten URLs were manually collected, resulting in an initial list of 1500 URLs representing the highest-ranking web pages for each search term. Duplicate URLs and keywords were excluded, resulting in a final count of 1427 URLs. A small number of links were deemed invalid, such as combined or PDF files, and were excluded from the analysis.

| No. | Feature | Acronym | Definition |
|-----|---------|---------|------------|
| 1 | **Amount of text** | total_words | Count the number of characters in paragraph and titles (<p> and <h> elements) |
| 2 | **H1 count of titles** | h1_num | Count the number of H1 titles on page |
| 3 | **H1 length** | h1_len | Count the average length of H1 titles on page |
| 4 | **H2 count of titles** | h2_num | Count the number of H2 titles on page |
| 5 | **H2 length** | h2_len | Count the average length of H2 titles on page |

| 6 | **H3 count of titles** | h3_num | Count the number of H3 titles on page |
|---|---|---|---|
| 7 | **H3 length** | h3_len | Count the average length of H3 titles on page |
| 8 | **Header total** | header_total | Count of all the headers on page |
| 9 | **Image count** | img_count | Count the number of images |
| 10 | **Internal links count** | internalLinks | Count the number of internal links (internal = linking to a page in the same domain) |
| 11 | **External links count** | externalLinks | Count the number of external links (external = linking to a page in the different domain) |
| 12 | **Total links count** | total_link | Count the number of total links (total links = internal links + external links) |
| 13 | **Keyword count H1** | h1_kcount | Count how many times the keyword mentioned in all the H1s |
| 14 | **Keyword count H2** | h2_kcount | Count how many times the keyword mentioned in all the H2s |
| 15 | **Keyword count** | h3_kcount | Count how many times the keyword |

| | | | |
|---|---|---|---|
| | **H3** | | mentioned in all the H3s |
| 16 | **Keyword count p** | p_kcount | Count how many times the keyword mentioned in all the paragraphs |
| 17 | **Keyword in anchor text** | a_kcount | 0 if keyword not in anchor text of any link, 1 if keyword in anchor text of any link |
| 18 | **Keyword in footer** | footer_kcount | 0 if keyword not in footer, 1 if keyword in footer |
| 19 | **Keyword in URL** | link_kcount | 0 if keyword not in URL, 1 if keyword in URL |
| 20 | **Keywords in image alt** | imalt_kcount | Count the number of times keyword mentioned in alt tag of images |
| 21 | **Meta desc length** | meta_desc_len | Count the length of the meta description. If no meta description, length = 0 |
| 22 | **Meta keywords count** | meta_kcount | Count the number of meta keywords used |
| 23 | **Page title used** | ti_used | 0 if no page title tag used, 1 if page title tag used |

*Table 2: Features Extracted from Webpages*

Overall, the data collection process for this scientific research entails a series of steps, including the identification of relevant keywords related to various aspects of the education industry, followed by a Google search to identify the top 10 webpages with the highest ranking for each keyword. Subsequently, a careful selection of the most relevant webpages is conducted, from which relevant data is extracted using 23 features. The extracted data is then subjected to a rigorous cleaning and preprocessing procedure, which involves removing irrelevant information, standardizing data across different webpages, and transforming text data into a format suitable for machine learning algorithms.

## IV. MODEL DEVELOPMENT

### 1. Theory

### 1.1. Search Engine Optimization (SEO)

According to Combe (2015), Search Engine Optimization (SEO) is an internet marketing technique that aims to enhance a website's visibility and attract valuable traffic from search engines. The main objective of SEO is to achieve a higher ranking on search engine results pages (SERPs) for specific keywords and phrases [34]. SEO involves employing various strategies and methods to enhance a website's relevance, authority, and reliability. Through optimizing a website's content, structure, and performance, SEO aims to facilitate search engine crawling and indexing while providing users with the most pertinent and beneficial search results for their queries.

In this paper, we will discuss the various types of SEO that are used to optimize web pages for search engines. There are several types of SEO, including:

On-page SEO pertains to the optimization of individual web pages in order to enhance their visibility and relevance on search engines. This involves optimizing elements like title tags, meta descriptions, header tags, and content to align with relevant keywords. On-page SEO holds significance for two primary reasons: it aids search engines in comprehending the page's content and assists users in finding the information they seek [35].

Off-page SEO involves employing strategies to enhance a website's reputation and authority through the acquisition of high-quality backlinks from reputable websites, social media sharing, and directory listings. The significance of off-page SEO lies in the fact that search engines utilize backlinks as a metric to gauge the popularity and relevance of a webpage [36].

Technical SEO encompasses the optimization of a website's technical elements, including page speed, mobile responsiveness, URL structure, and sitemap optimization. The significance of technical SEO lies in its ability to facilitate efficient crawling and indexing of the website by search engines, while also enhancing the overall user experience [37].

Local SEO focuses on optimizing a website for local search by enhancing Google My Business profiles, local directory listings, and creating location-specific content. This type of SEO is particularly important for businesses with physical locations or those targeting specific geographic areas. It enables these businesses to improve their visibility and reach within their local community, ultimately driving more relevant traffic and potential customers to their doorstep [38].

E-commerce SEO involves optimizing a website specifically for e-commerce searches, focusing on elements such as product descriptions, pricing, and the optimization of shopping carts and checkout processes. This type of SEO is crucial for online retailers who sell products through their website. By implementing effective e-commerce SEO strategies, online retailers can improve their visibility in search results, attract targeted traffic, and increase the likelihood of conversions and sales [39].

International SEO involves optimizing a website for international searches by employing techniques such as geotargeting, hreflang tags, and content localization. This type of SEO is particularly significant for businesses that cater to customers in different countries and languages. By implementing international SEO strategies, businesses can enhance their visibility and relevance in global search results, effectively targeting and engaging with their international audience. This, in turn, facilitates increased traffic, improved user experience, and higher chances of conversions and business growth on a global scale [40].

SEO can be applied to any platform that has a search engine, including search engines themselves (like Google or Bing), social media platforms (like Facebook or Instagram), e-commerce platforms (like Amazon or eBay), and more. Essentially, any platform that relies on users searching for information can benefit from SEO to improve the visibility and relevance of its content.

## 1.2. Machine Learning

Machine learning is a field of artificial intelligence that focuses on refining algorithms through the use of empirical data [41]. The primary objective of machine learning is to enable machines to acquire knowledge from data, allowing them to make predictions, identify patterns, and uncover underlying structures [42]. Supervised learning, a widely used approach within machine learning, involves working with labeled data, where the values of the dependent variable are known [42]. In the context of web page ranking on search engines, supervised learning techniques can be applied to develop models that predict the position of a webpage in search engine rankings, considering a wide range of features and factors [43]. The specific methodologies and algorithms employed in supervised learning for web page ranking can vary, but the core concept revolves around using a labeled dataset where each instance corresponds to a webpage and its associated position in search engine rankings.

*1.2.1. The ranking problem*

In the context of web page ranking, the ranking problem involves the task of arranging web pages in order of relevance to a user's query, aiming to provide accurate and useful search results, while the evaluation of ranking models is often done using metrics like NDCG (Normalized Discounted Cumulative Gain) that consider the relevance scores and positions of webpages in the ranked list [43]. NDCG takes into account both the relevance scores assigned to individual webpages and the position of each web page in the ranked list. The calculation of NDCG involves two main components: DCG (Discounted Cumulative Gain) and IDCG (Ideal DCG). Mathematically, NDCG can be expressed as:

$$nDCG_k = \frac{DCG_k}{IDCG_k}$$

where $nDCG_k$ represents the NDCG score at position $k$ in the ranked list. $DCG_k$ (Discounted Cumulative Gain) calculates the accumulated relevance scores of the top (k) webpages, with decreasing weights assigned to each position. Wheres $IDCG_k$ (Ideal DCG) denotes the maximum achievable DCG score at position (k) when the webpages are perfectly ordered based on their relevance.

The formula for DCG is as follows, where $rel_i$ represents the relevance score of the $i^{th}$ web page in the ranked list.

$$DCG_k = \sum_{i=1}^{k} \frac{rel_i}{log_2(i + 1)}$$

The formula for IDCG (Ideal DCG) is calculated in a similar manner to DCG (Discounted Cumulative Gain), but assumes a perfect ordering of web pages based on their relevance scores. By dividing the DCG by the IDCG, NDCG provides a normalized measure of the

21

ranking performance, ranging from 0 to 1, where a higher NDCG score indicates a better-ranked list with more relevant web pages positioned at the top [44].

### 1.2.2. Gradient boosting

Determining the relevance of documents or web pages to a given query is a crucial aspect of information retrieval systems. In the past, humans used to assess relevance scores based on their subjective judgments. However, thanks to the advancements in machine learning and natural language processing, we now have models that can automatically determine these scores [45].

In the field of web page ranking, gradient boosting emerges as a widely acclaimed framework, particularly due to its ability to determine relevance scores even for data that lacks pre-assigned relevance scores. This capability is invaluable when dealing with large datasets where manual labeling of relevance scores for each data point is impractical or time-consuming. By leveraging the power of gradient boosting, web page ranking models can autonomously learn and infer relevance scores based on patterns and relationships present in the available data [46].

Gradient boosting frameworks are widely employed as part of the supervised learning approach in machine learning. These frameworks aim to construct a strong predictive model by sequentially combining multiple weak learners, typically decision trees.The fundamental idea behind gradient boosting is to iteratively minimize the overall prediction error by focusing on the errors made by the previous iterations. This is achieved by optimizing a loss function through a gradient descent algorithm, where each subsequent weak learner is built to address the shortcomings of the previous ones. The final prediction is obtained by aggregating the predictions of all the weak learners, weighted by their respective contribution to the overall model [43].

Mathematically, the gradient boosting algorithm can be represented as follows:

$$F(x) = \sum_{m=1}^{M} \gamma_m h_m(x)$$

where $F(x)$ represents the final prediction for input $x$, $M$ is the total number of weak learners, $\gamma_m$ denotes the weight assigned to the $m^{th}$ weak learner, and $h_m(x)$ corresponds to the prediction made by the $m^{th}$ weak learner.

Gradient boosting algorithms, such as LightGBM Ranker and XGBoost Ranker, have gained popularity for their ability to capture intricate relationships between webpage features and search engine rankings [47]. These frameworks excel in handling large-scale datasets efficiently and offer advanced features such as parallelization options, regularization techniques, and early stopping criteria [48].

## 2. Model selection

For this research, we have chosen to utilize two of the most efficient models for learning-to-rank (LTR) tasks: LightGBM Ranker and XGBoost Ranker. These models are renowned for their speed, scalability, and flexibility in handling large and high-dimensional datasets, making them ideal for predicting webpage rankings in the healthcare industry.

LightGBM Ranker and XGBoost Ranker have been shown to outperform traditional methods such as linear regression, logistic regression, and Support Vector Machines (SVM) in various LTR tasks [49] [50]. One of the key reasons for their superior performance is their ability to efficiently handle large-scale datasets. For instance, in a comparative study on LTR for web search, LightGBM Ranker was found to be significantly faster than other popular methods such as LambdaMART and RankSVM [51]. This efficiency is particularly beneficial for our study, which involves analyzing extensive datasets derived from web scraping healthcare-related keywords.

Both LightGBM Ranker and XGBoost Ranker offer a range of objective functions and evaluation metrics, providing greater flexibility in modeling and evaluation. LightGBM Ranker supports various objective functions, including pairwise, ndcg, and lambdarank, while XGBoost Ranker offers objective functions such as rank [52] [53]. This flexibility allows for the selection of the most appropriate objective function to optimize the model's performance for specific ranking tasks.

Moreover, these models have demonstrated state-of-the-art performance in various LTR applications. For example, in a study on LTR for product search, XGBoost Ranker outperformed other popular methods such as LambdaMART and RankSVM in terms of Normalized Discounted Cumulative Gain (NDCG) [54]. Given these capabilities, LightGBM Ranker and XGBoost Ranker are well-suited for our objective of predicting the ranking of web pages in the healthcare industry and identifying factors that influence these rankings.

**Implementation of LightGBM Ranker and XGBoost Ranker**

Although LightGBM Ranker and XGBoost Ranker differ in their implementation details, both frameworks employ similar gradient boosting techniques to optimize the loss function and predict the rank of new inputs.

LightGBM Ranker utilizes a leaf-wise tree growth strategy, which speeds up the training process and supports a wide range of objective functions and evaluation metrics [52]. This strategy allows the model to grow tree leaves more efficiently, leading to faster computation and reduced memory usage. The ability to handle large datasets efficiently makes LightGBM particularly advantageous for our study.

On the other hand, XGBoost Ranker combines gradient boosting with regularization techniques to minimize the loss function and handle sparse data efficiently [53]. XGBoost's regularization capabilities prevent overfitting, ensuring that the model generalizes well to

new, unseen data. This is crucial for maintaining the accuracy and reliability of the predicted rankings in a real-world healthcare context.

**Gradient Boosting Mechanism**

Gradient boosting iteratively adds weak learners to the model, typically decision trees, to optimize the loss function. The process involves the following steps:

1. **Initialization:** The model starts with an initial prediction.
2. **Residual Calculation:** For each iteration, the residual (difference between the actual and predicted values) is calculated.
3. **Fitting Weak Learners:** A new decision tree is fitted to the residuals. The goal is to approximate the negative gradient of the loss function concerning the current model's predictions.
4. **Update:** The model's predictions are updated by adding the newly fitted decision tree, scaled by a learning rate.
5. **Iteration:** Steps 2-4 are repeated for a predefined number of iterations or until convergence.

This iterative process allows the model to minimize the loss function by learning optimal weights for the associated features. As new decision trees are added, the model progressively improves its predictions, thereby enhancing the ranking accuracy.

Once the model is trained, it can predict the rank of new input by feeding the input features into the ensemble and aggregating the predictions from all the decision trees. The output is a ranking score representing the predicted relevance of the search results [55].

By deploying LightGBM Ranker and XGBoost Ranker, this study aims to achieve accurate and efficient predictions of webpage rankings in the healthcare industry. The insights gained from these models will help identify the key factors influencing webpage rankings and inform strategies to enhance the online visibility of healthcare information.

### 3. Model evaluation

**Training, Validation, and Testing**

The model evaluation process was meticulously designed to ensure the robustness and reliability of the machine learning models. The dataset was divided into three subsets: training, validation, and testing. The training subset was used to train the models, while the validation subset was employed to tune hyperparameters and avoid overfitting. The testing subset was reserved to evaluate the generalization performance of the models, providing a more accurate assessment of their prediction power.

**Handling Groups in Learning-to-Rank**

A unique aspect of learning-to-rank (LTR) problems is the presence of "groups" within the data. In this context, each group corresponds to a set of search results for a specific keyword. The ranker models utilized in this study required not only the training, validation, and testing sets but also the information on the number of samples within each group. This grouping is essential as the ranking of search results is highly dependent on the search query context and the relevance of results to the user's information needs.

**Evaluation Metric: Normalized Discounted Cumulative Gain (NDCG)**

To evaluate the performance of the models, Normalized Discounted Cumulative Gain (NDCG) was used as the primary ranking metric. NDCG is a widely adopted metric in the field of LTR, reflecting the quality of the ranked results. The fundamental principle behind NDCG is to discount the relevance score of each result based on its position in the ranking and normalize these discounted scores by dividing them by the ideal discounted score, which represents the best possible ranking for a given set of results [50]. The formula for NDCG is as follows:

$$\text{NDCG} = \frac{\text{DCG}}{\text{IDCG}}$$

where DCG (Discounted Cumulative Gain) is calculated as:

$$\text{DCG} = \sum_{i=1}^{p} \frac{2^{rel_i} - 1}{\log_2(i+1)}$$

and IDCG (Ideal Discounted Cumulative Gain) is the DCG score of the ideally ranked documents.

**Model Performance**

The performance of the two models on the testing set was evaluated using the NDCG metric. The results are summarized below:

| Model | NDCG Score on the testing set |
|---|---|
| XGBoost | 0.877 |
| Light GBM | 0.856 |

*Figure 1: Average NDCG Scores of Models.*

Based on these results, XGBoost was selected for further analysis and feature importance evaluation due to its superior performance. The higher NDCG score indicates that XGBoost provided more relevant rankings of search results compared to LightGBM. One possible reason for XGBoost's outperformance is its ability to handle sparse data more efficiently than LightGBM. The dataset in this study contained multiple zero values due to the absence of certain features in specific webpages, rendering it sparse. XGBoost utilizes a "sparse-aware" gradient boosting technique, which capitalizes on the sparsity structure of the data. This technique updates only the non-zero gradients and Hessians during the gradient boosting process, significantly reducing the computational cost and memory usage of the algorithm. Moreover, XGBoost employs L1 regularization to encourage sparsity in the learned model. L1 regularization forces many feature weights to be set to zero, effectively

selecting only the most important features and reducing the model's complexity. This ability to handle sparse data efficiently and regularize the model contributes to XGBoost's superior performance in this study [56].

In contrast, LightGBM does not have a built-in mechanism for handling sparse data as effectively. It treats all features as dense and does not leverage the sparsity structure of the data. Consequently, LightGBM may require more memory and computational resources to process sparse data, particularly when the number of features is high [57]. This limitation can hinder its performance relative to XGBoost in scenarios involving sparse datasets.

By leveraging XGBoost's capabilities, this study was able to achieve a more accurate and efficient prediction of webpage rankings in the healthcare industry. The insights derived from the model's performance and feature importance analysis can inform strategies to enhance the online visibility of healthcare information, ultimately benefiting users seeking reliable health-related content.

## 4. Features importance analysis

The primary goal of our study is to identify the most critical factors that influence webpage rankings and to enhance these features within the specific domain of healthcare. To achieve this, we utilized the XGBoost model and assessed feature importance using two methods: the F-score and SHAP (SHapley Additive exPlanations) values.

**F-score Method**

The F-score is a straightforward method that represents the number of times each feature is used to split the data across all trees in the model. This metric provides an indication of the feature's significance based on its frequency of use in the decision-making process. Although useful, the F-score does not account for the context of each split or the feature's overall contribution to the model's predictive performance.

**SHAP Values**

SHAP values offer a more nuanced approach to feature importance analysis. This technique interprets the output of machine learning models by attributing the importance of each feature to the final prediction. Based on the concept of Shapley values from cooperative game theory, SHAP values calculate the contribution of each feature by estimating the expected change in the model's output when the feature is included or excluded.

One significant advantage of SHAP values is that they provide both global and local explanations of the model's behavior. Global explanations offer an overall assessment of the importance of each feature across the entire dataset, while local explanations allow for an understanding of how a specific feature contributes to a prediction for an individual instance. This dual perspective is particularly valuable in interpreting complex models and identifying specific areas for improvement [56].

**Comparison of SHAP and F-score**

In a study published in the Journal of Healthcare Informatics Research, the authors compared SHAP values and F-score for feature importance analysis in a binary classification problem using XGBoost. They found that SHAP values provided more accurate and interpretable feature importance rankings compared to the F-score [58]. Similarly, a study in the Journal of Machine Learning Research evaluated several feature importance methods, including SHAP and F-score, in a regression problem using XGBoost and concluded that SHAP was the most effective method for identifying important features [59]. Additionally, research published in the Journal of Applied Sciences compared the performance of SHAP and F-score for feature importance analysis in a prediction problem using XGBoost and determined that SHAP provided more accurate and reliable results [60].
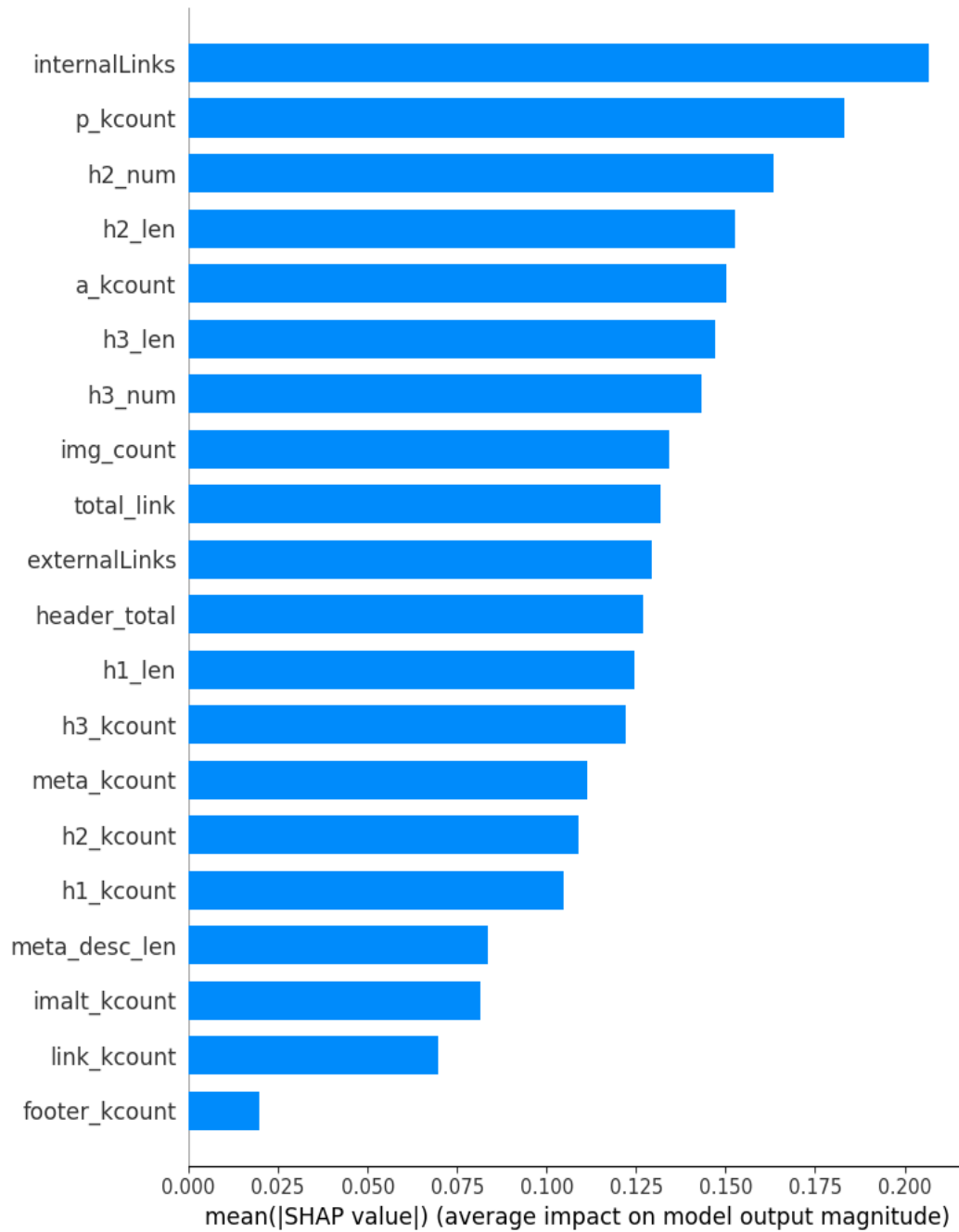
## 4.1. Application of SHAP in Our Study

Given the advantages and superior performance of SHAP values demonstrated in previous research, our primary method for evaluating the impact of features on webpage ranking is SHAP.

The SHAP beeswarm summary plot below displays the distribution of SHAP values for each feature. The plot consists of a horizontal axis representing the SHAP value for each instance in the dataset and a vertical axis representing the feature. Each dot in the plot depicts an instance, and its position on the x-axis corresponds to its SHAP value for a particular feature. Features with positive SHAP values indicate that the feature contributes to a higher output from the model, while features with negative SHAP values suggest that the feature contributes to a lower output. The features are sorted in descending order of their mean SHAP value, with the most important feature at the top.

This visualization provides a comprehensive overview of the relative importance of each feature and its impact on the model's predictions, enabling us to identify key areas to focus on for enhancing webpage rankings in the healthcare domain.

The SHAP beeswarm summary plot below provides a detailed view of the impact of various features on webpage ranking. Each dot represents an instance, and its position along the horizontal axis corresponds to the SHAP value for that feature. The color gradient from blue to red indicates the feature value from low to high. Features with positive SHAP values contribute to a higher model output, indicating a higher rank, while features with negative SHAP values contribute to a lower model output. The features are sorted in descending order of their mean SHAP value, with the most influential feature at the top.
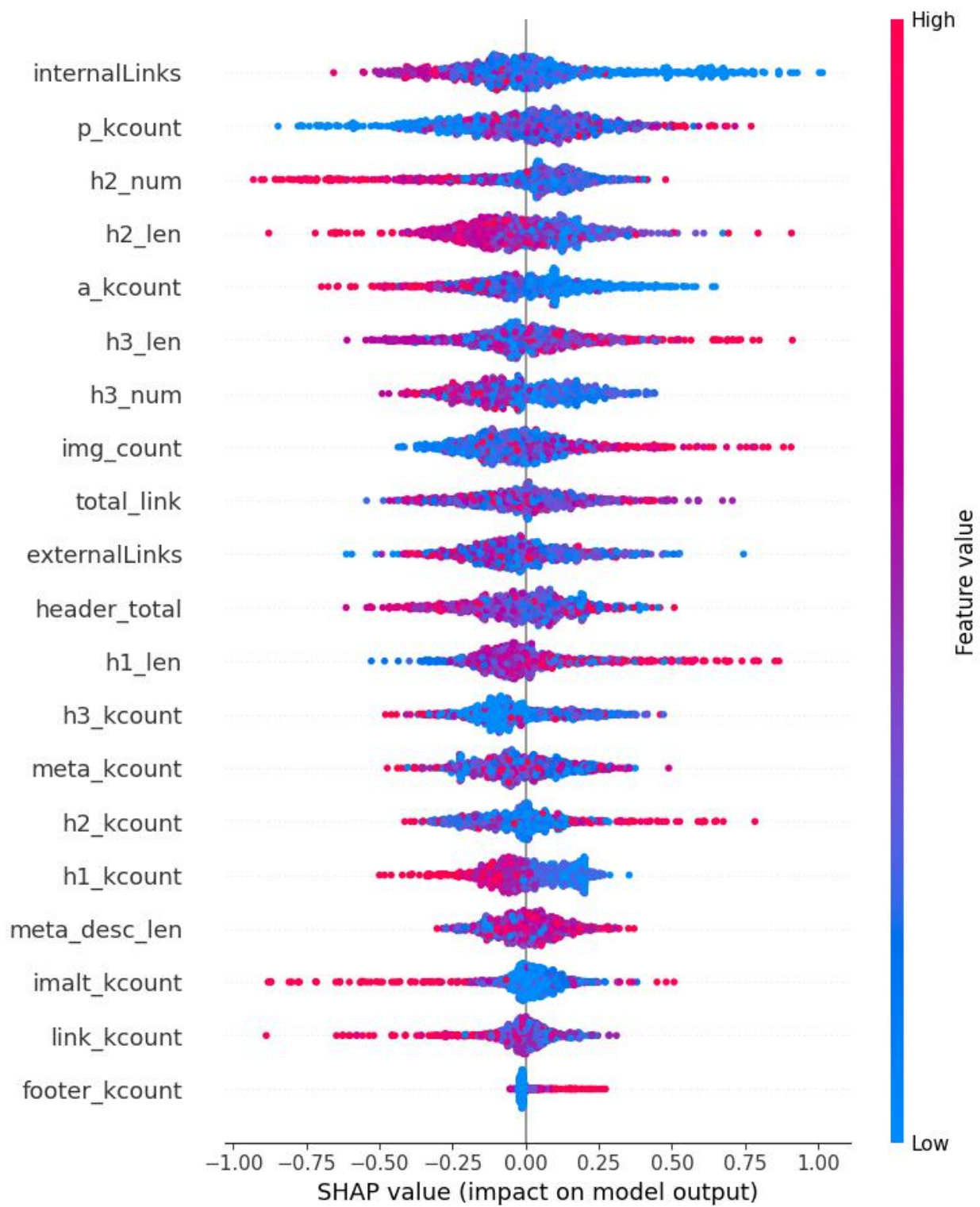
*Figure 2: SHAP beeswarm summary plot measured for XGBoost*

-

**Key Insights from the SHAP Beeswarm Plot:**

- **Internal Links (internalLinks):** The number of internal links is the most significant feature influencing webpage ranking. High internal link values (indicated by red dots) are associated with higher SHAP values, suggesting that a well-structured network of internal links can significantly boost a webpage's ranking. This highlights the importance of a well-organized internal linking strategy, which can improve site navigation and enhance the SEO value of the content. However, it is important to avoid excessive internal linking, which can dilute content relevance and make it harder for search engines to crawl and understand the page.

- **Paragraph Keyword Count (p_kcount):** The frequency of keywords within paragraphs is also a crucial factor. High keyword density (red dots) within the content correlates with higher SHAP values, indicating a positive impact on rankings. This underscores the need for relevant and keyword-rich content that addresses the users' queries effectively.

- **Number and Length of H2 Tags (h2_num, h2_len):** Both the number and length of H2 tags are important. Higher numbers of H2 tags (red dots) are associated with positive SHAP values, emphasizing the role of well-structured headings in content organization. Shorter H2 tags (blue dots for lower length values) are generally more favorable, indicating that concise and relevant headings can improve readability and SEO performance.

- **Anchor Keyword Count (a_kcount):** The frequency of keywords within anchor tags (<a> tags) also significantly influences rankings. Higher counts (red dots) are positively associated with SHAP values, suggesting that relevant and keyword-rich anchor text enhances the SEO value of hyperlinks.

- **Number and Length of H3 Tags (h3_num, h3_len):** Similar to H2 tags, the number and length of H3 tags contribute to content structuring. Proper use of H3 tags helps further refine the content hierarchy, aiding both user experience and search engine understanding.

- **Image Count (img_count):** The number of images on a webpage affects rankings. A moderate number of well-optimized images (represented by red and blue dots evenly spread) can enhance user engagement without negatively impacting load times. However, too many images can slow down the page, adversely affecting rankings.

- **Total Links (total_link):** The total number of links, both internal and external, impacts the ranking. While links contribute to the page's connectivity and authority, an optimal balance is necessary to avoid overwhelming the content with too many links.

- **External Links:** The presence of external links must be managed carefully. Excessive external links (blue dots) can have a negative impact on rankings, as they may divert traffic and reduce the page's authority.

- **Header Total (header_total):** The cumulative length of all header tags (h1, h2, h3) is another important factor. Properly structured and concise headers improve content readability and SEO.

- **Length of H1 Tag (h1_len):** The length of the H1 tag is critical. Shorter, more descriptive H1 tags (blue dots) are associated with higher SHAP values, indicating their positive impact on ranking.

- **Meta Description Length (meta_desc_len):** The length of the meta description also influences rankings. A well-crafted meta description can improve click-through rates from search engine results pages (SERPs).

- **Miscellaneous Factors:** Other factors such as the keyword count in headers (h1_kcount, h2_kcount, h3_kcount), meta keywords (meta_kcount), and footer keyword count (footer_kcount) also play roles in influencing rankings, albeit to a lesser extent.

## 4.2. In-Depth Analysis of SEO Feature Characteristics
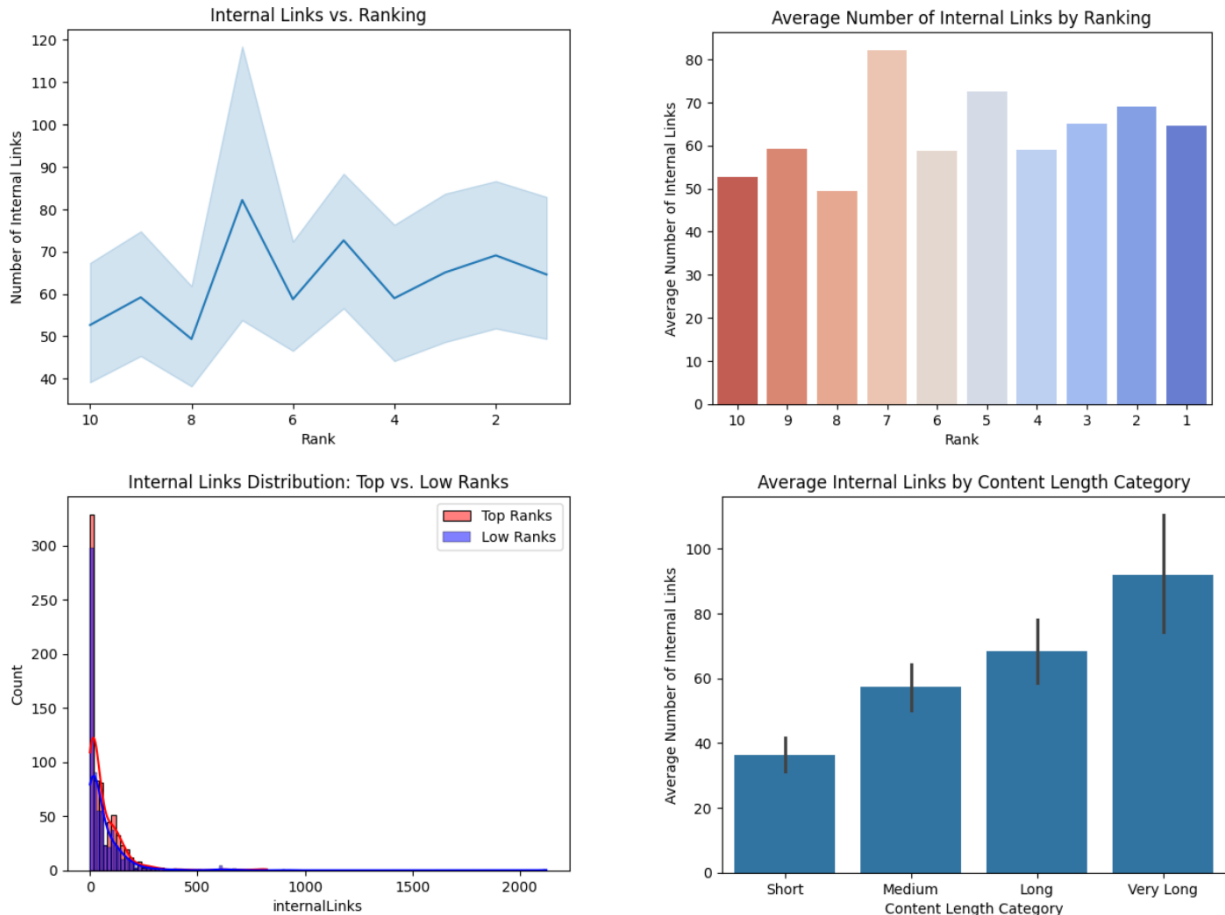
## Internal Links



*Figure 3: In-Deep Analysis of Internal Links*

- **Top Ranking Pages**: Ranks 1-2 show a consistent average number of internal links (60-70), with rank 7 showing the highest average (80).
- **Content Length**: Longer content is associated with more internal links, with very long content averaging about 90 internal links.
- **Peak Rank 7**: Notably higher internal links (~110) compared to other ranks.
- **Distribution**: Top ranks show broader internal link distribution, indicating a strategic use of internal links across different content lengths and types.

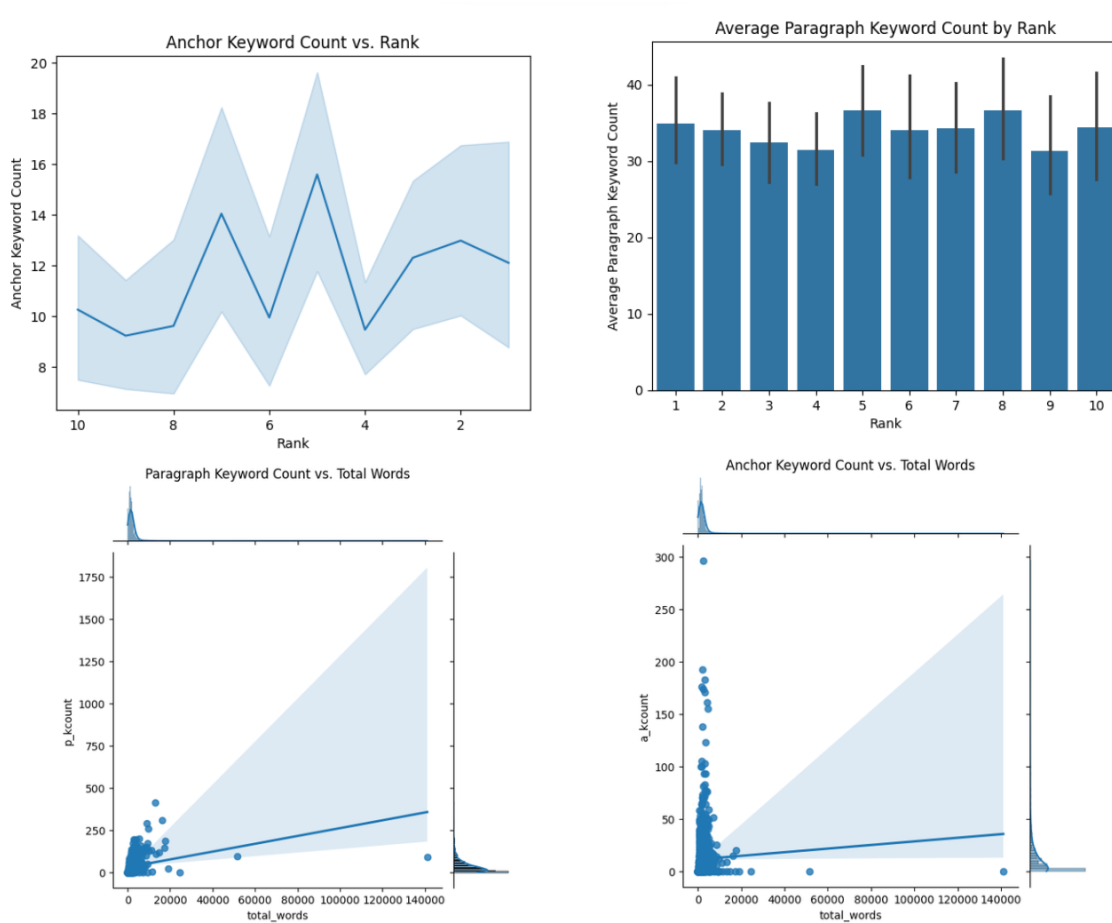**Analysis of Paragraph Keyword Count and Anchor Keyword Count**



*Figure 4: In-Deep Analysis of Paragraph Keyword Count and Anchor Keyword Count*

**Anchor Keyword Count**:

- Top-ranked pages (1-2) typically have 12-14 anchor keywords.
- Rank 6 shows a peak with approximately 18 anchor keywords, indicating variability in the optimal number of anchor keywords.

**Paragraph Keyword Count**:

- Top-ranked pages (1-2) have an average paragraph keyword count of 30-35.
- Slight increase at ranks 5-6 with averages around 35-40.

**Total Words Relationship**:

- **Paragraph Keyword Count**: Positive correlation with total words, but keyword density decreases with more words.
- **Anchor Keyword Count**: Slight positive correlation with total words, with higher keyword counts more frequent in shorter content.

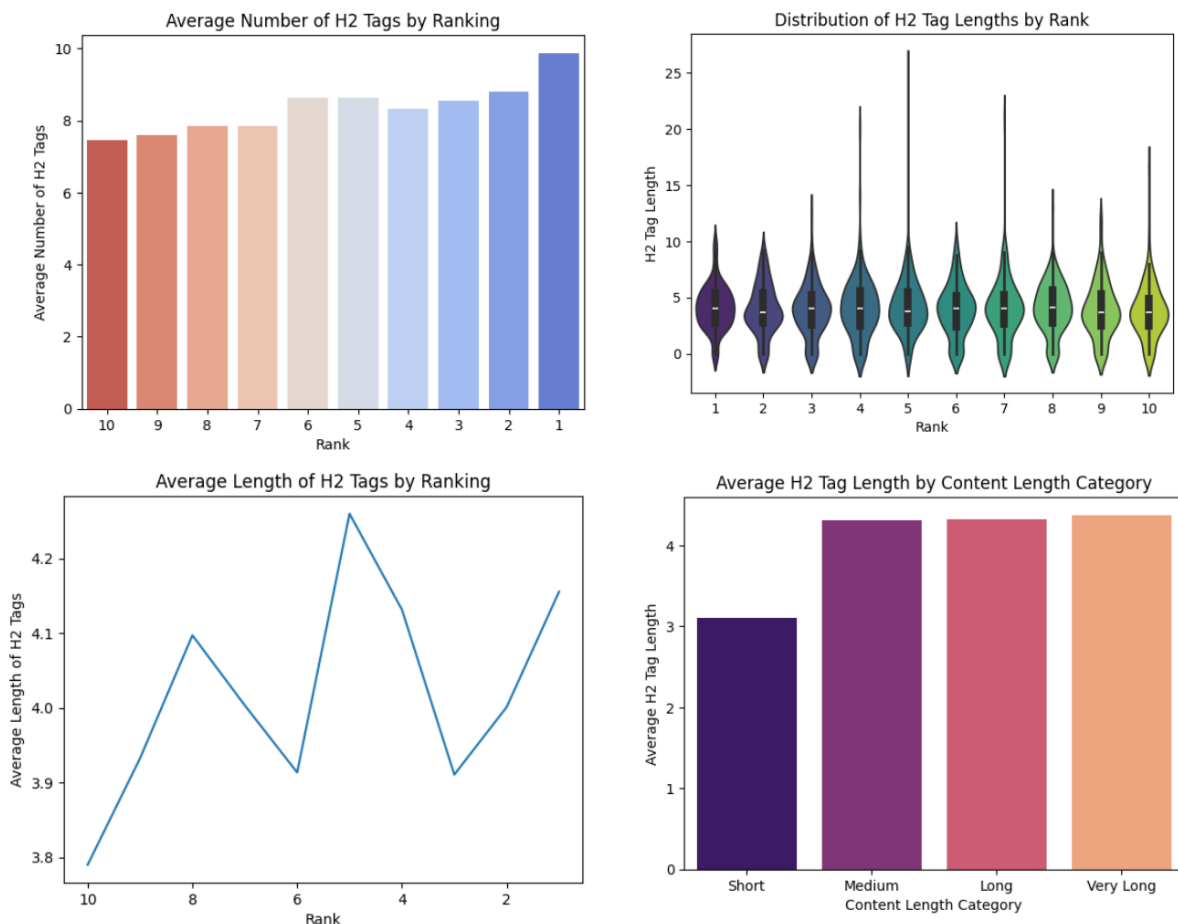**Number of H2 tags and Length of H2 tags**



*Figure 5: In-Analysis of Number of H2 tags and Length of H2 tags*

**Number of H2 Tags**:

- **Top-ranked pages (1-2)** have the highest average number of H2 tags (around 10).

- **Lower-ranked pages (10)** have fewer H2 tags on average (around 7).

**Length of H2 Tags**:

- **Top-ranked pages (1-2)** tend to have concise H2 tags with an average length of 4.0-4.2 characters.
- There is more variability in H2 tag lengths among lower ranks.

**Content Length Dependency**:

- **Short Content**: Slightly shorter H2 tags (average 3.8 characters).
- **Medium to Very Long Content**: Consistent average H2 tag length around 4.1 characters.

**General Insights**:

- Higher-ranked pages use more H2 tags, indicating better content structuring.
- Concise H2 tags are common among top-ranked pages, suggesting that clear and succinct headings contribute positively to rankings.
- The number and length of H2 tags should be tailored to content length, ensuring appropriate structuring and readability.
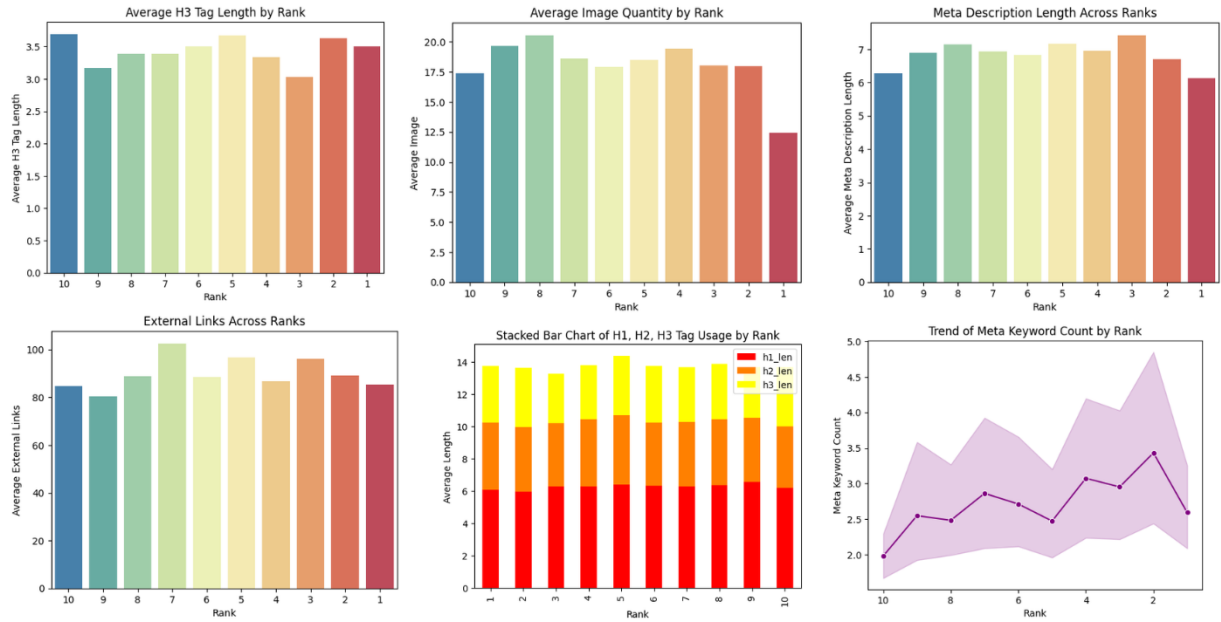
**Analysis of Another features**

Figure 6: Analysis of H3 Tag Length, Image Quantity, Meta Description Length, External Links, H1, H2, H3 Tag Usage, Meta Keyword Count

**H3 Tag Length**:

- **Top-ranked pages (1-3)** have slightly longer H3 tags (around 3.5 characters).
- **Lower-ranked pages (10)** have shorter H3 tags (around 3.2 characters).

**Image Quantity**:

- **Top-ranked pages (1-4)** have fewer images on average (13-18 images).
- **Rank 7** has the highest average number of images (21 images).

**Meta Description Length**:

- **Top-ranked pages (1-4)** have slightly longer meta descriptions (6.7-7 characters).
- **Lower-ranked pages (10)** have shorter meta descriptions (6.3 characters).

**External Links**:

- **Top-ranked pages (1)** have fewer external links (77 on average).

- **Rank 7** has the highest number of external links (96 on average).

**H1, H2, H3 Tag Usage**:

- **Consistent Length**: Across ranks, H1 tags average around 5-6 characters, H2 tags around 4 characters, and H3 tags around 3-4 characters.
- **General Insight**: Maintaining consistent tag lengths might be beneficial for SEO.

**Meta Keyword Count**:

- **Top-ranked pages (2-3)** have higher meta keyword counts (4-4.5).
- **Lowest-ranked pages (10)** have fewer meta keywords (2.2 on average).

## V. DISCUSSION

### 1. Practical Implications to Marketers for SEO in the Healthcare Industry

*The findings from our study highlight several practical implications for healthcare marketers aiming to enhance their SEO efforts. The analysis particularly emphasizes the importance of internal links, paragraph keyword count, H2 tags, anchor keyword count, and H3 tags as the most impactful features. Here, we delve into these features in a detailed and practical manner, providing actionable strategies supported by examples and case studies to illustrate their application.*

**Internal Links**

Internal links emerged as the most significant feature influencing webpage ranking. A well-structured internal linking strategy can significantly boost a webpage's SEO performance by improving site navigation and enhancing the SEO value of the content. For instance, Mayo Clinic's website exemplifies effective internal linking, where each health article links to related conditions, treatments, and patient stories, creating a comprehensive internal web.

To implement this, healthcare marketers should develop a strategic internal linking plan that ensures all relevant pages within the site are interconnected. This can be achieved by using tools like Yoast SEO for WordPress, which automates internal linking suggestions, ensuring consistency and saving time. Additionally, marketers should prioritize quality over quantity by auditing current internal links and removing or updating those that do not provide additional value or context. Moreover, linking new content from high-authority pages can help distribute link equity throughout the site, enhancing the overall SEO impact.

**Paragraph Keyword Count**

The frequency of keywords within paragraphs is crucial for SEO performance. High keyword density within the content correlates with higher SHAP values, indicating a positive impact on rankings. Cleveland Clinic effectively uses tools like Ahrefs and SEMrush to conduct keyword research and integrate relevant keywords naturally within their content, thereby enhancing their SEO performance.

Healthcare marketers should conduct comprehensive keyword research using tools like Google Keyword Planner or Ahrefs to identify relevant keywords that align with their target audience's search behavior. These keywords should be integrated naturally into the content to avoid keyword stuffing, which can negatively impact readability and SEO. Focusing on long-tail keywords can attract more targeted traffic, while periodically updating content with new and relevant keywords can maintain its SEO value.

**Number and Length of H2 Tags (h2_num, h2_len)**

Both the number and length of H2 tags play a significant role in content organization and SEO performance. Higher numbers of H2 tags are associated with positive SHAP values, while concise H2 tags are generally more favorable. The NHS website effectively uses H2 tags to break content into sections, making it easier to navigate and understand.

To optimize H2 tags, healthcare marketers should break down their articles into clear, concise sections using H2 tags for each main point. This approach enhances readability and

helps search engines understand the content structure. It is essential to keep H2 tags concise and descriptive, ideally between 4-6 words, to ensure clarity and relevance. Additionally, including primary keywords in H2 tags can signal the relevance of the section to search engines, further enhancing SEO performance.

**Anchor Keyword Count (a_kcount)**
The frequency of keywords within anchor tags significantly influences webpage rankings. Relevant and keyword-rich anchor text enhances the SEO value of hyperlinks. Healthline optimizes anchor text by using descriptive and keyword-rich phrases, such as "symptoms of high blood pressure," rather than generic terms like "click here."

Healthcare marketers should optimize anchor text by ensuring it is descriptive and keyword-rich, accurately reflecting the linked content. Contextual relevance is crucial; anchor text should be placed within relevant content to improve user understanding and engagement. Regularly reviewing and updating anchor texts can maintain their relevance and effectiveness, while leveraging internal anchor texts can create a robust internal linking structure, distributing link equity effectively.

**Number and Length of H3 Tags (h3_num, h3_len)**

Similar to H2 tags, the number and length of H3 tags contribute to content structuring and SEO performance. Proper use of H3 tags helps further refine the content hierarchy, aiding both user experience and search engine understanding. Cleveland Clinic uses H3 tags to create subsections within H2 sections, improving readability and navigation.

Healthcare marketers should use H3 tags to break down content within H2 sections, creating manageable and understandable subsections. Keeping H3 tags concise and relevant, typically 3-4 words, ensures clarity and improves readability. Including relevant keywords in H3 tags can enhance their SEO value, while maintaining a logical and consistent use of H3 tags throughout the content can improve overall structure and user experience.

By focusing on these practical SEO strategies, healthcare marketers can significantly improve their website rankings and visibility. Emphasizing internal links, keyword-rich content, optimized headers, descriptive anchor texts, and well-structured content hierarchies can help healthcare websites provide valuable information to a broader audience, ultimately enhancing user experience and engagement. These strategies, supported by real-world examples and case studies, offer a roadmap for healthcare marketers to implement effective SEO practices and achieve better search engine rankings.

## 2. Limitations and Future Works

One limitation of our study is the focus on the healthcare industry. While this allowed for specific and nuanced insights into SEO for this sector, the findings may not necessarily generalize to other industries. Future research could explore the impact of these features on SEO in various other sectors, such as finance, e-commerce, or education. Additionally, our proprietary Python code and methodology are available in our software, facilitating data collection and effective industry-specific analysis, thereby yielding more accurate outcomes.

Furthermore, our study concentrated on a specific set of features identified through feature analysis. However, there may be other factors not included in our analysis that could also influence SEO. To provide a more comprehensive understanding of SEO, future research should explore additional features beyond those examined in our study. These might include elements such as website traffic, social media influence, and customer behavior on the website. Analysis of these additional features could potentially reveal new insights into how search engines evaluate and rank webpages, helping businesses and organizations to develop more effective SEO strategies.

In terms of technological advancements, our study did not fully utilize social media metrics and user engagement data, such as click-through rates and dwell time, which are critical components of SEO. Integrating these data points could offer a more holistic view of the

factors influencing webpage rankings. Future research should aim to include these variables to enhance the comprehensiveness of the analysis.

In the future, we plan to continue developing our SEO tool to include additional features beyond internal link suggestions. Specifically, we aim to incorporate capabilities for external link generation, optimal anchor text creation, title suggestions and edits, and automatic identification of SEO errors in content pages using the Chat GPT API. Our goal is to develop a comprehensive product that can benefit not only the healthcare industry but also other sectors seeking to improve their search engine rankings.

## VI.    REAL IMPLEMENTATION WITH EVERYDAY HEALTH WEBSITE

*In this section, we detail the application of our SEO research and methodologies to the website of Everyday Health, a prominent online resource providing health and wellness information. Everyday Health offers a wide range of articles, tools, and resources aimed at helping individuals manage their health and wellness. As a leading digital health company, it is essential for Everyday Health to maintain high search engine rankings to reach a broader audience and achieve its marketing objectives. This study involves crawling the webpages of Everyday Health, analyzing 23 key SEO features, and providing actionable recommendations to enhance their webpage rankings.*

### 1.  Everyday Health Introduction

Everyday Health, accessible at www.everydayhealth.com, is dedicated to providing authoritative health information, news, and tools for managing health conditions and wellness. The website covers a diverse range of health topics including chronic conditions, mental health, diet, fitness, and general wellness. Each webpage is designed to be informative and user-friendly, featuring articles, videos, infographics, and interactive tools to engage users and provide valuable health insights.

The content strategy of Everyday Health emphasizes high-quality, evidence-based information that is accessible to a wide audience. The website's structure includes comprehensive health guides, symptom checkers, and personalized health recommendations, all aimed at empowering users with the knowledge to make informed health decisions. Given the competitive nature of the digital health space, optimizing these webpages for search engines is critical to maintaining and improving visibility and user engagement.

## 2. Data Collection

To identify areas for improvement, we conducted a comprehensive crawl of the Everyday Health website, collecting data on 23 SEO-related features from 126 of their articles and webpages by using BeautifulSoup via Python. These features include internal and external links, keyword usage in headers and paragraphs, meta description length, image count, and more. The dataset created from this crawl was then analyzed and compared with the top 10 ranking insights gained from our previous research.

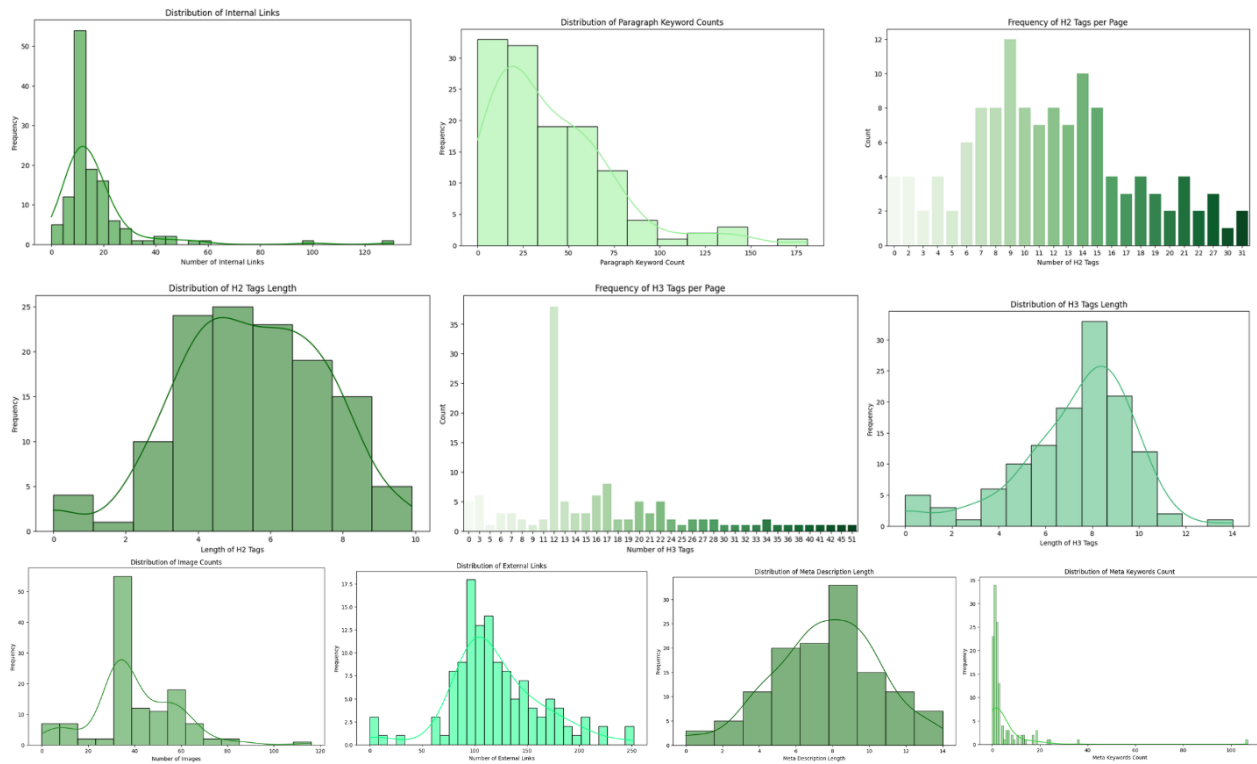## 3. Visualization and Comparison with Top Ranking Insights



*Figure 7: Distribution of SEO features in Everyday Health website*

## Distribution of Internal Links

The distribution of internal links on Everyday Health's webpages shows a right-skewed pattern, with most pages having fewer than 20 internal links, and a few pages having significantly more, up to around 120. In contrast, top-ranking pages (ranks 1-2) consistently have 60-70 internal links, with rank 7 showing a peak at approximately 110 internal links. This indicates that Everyday Health should aim to increase the number of internal links, targeting an average of 60-70, to improve SEO performance.

## Distribution of Paragraph Keyword Counts

Paragraph keyword counts exhibit a right-skewed distribution, with most pages having fewer than 50 keywords. Some pages have up to 175 keywords, but these are outliers. Top-

ranking pages (ranks 1-2) typically have 30-35 keywords per paragraph, with a slight increase at ranks 5-6 to 35-40. Everyday Health should increase keyword density within paragraphs to align with these benchmarks.

**Frequency of H2 Tags per Page**

The frequency of H2 tags per page shows a normal distribution, with most pages containing between 6 to 14 H2 tags, peaking around 10. Top-ranking pages (ranks 1-2) typically have around 10 H2 tags. Everyday Health is performing well in this aspect but should ensure all pages consistently use a high number of H2 tags for optimal SEO.

**Distribution of H2 Tags Length**

The length of H2 tags on Everyday Health pages varies, with an average length of around 4 to 6 words. This aligns with top-ranking pages, where concise H2 tags (around 4.0-4.2 characters) are preferred. Maintaining or slightly reducing the length of H2 tags will enhance clarity and SEO performance.

**Frequency of H3 Tags per Page**

The number of H3 tags per page varies widely, with some pages having no H3 tags and others having up to 50. Top-ranking pages effectively use H3 tags to structure content. Everyday Health should standardize the use of H3 tags across their pages to improve content organization and SEO.

**Distribution of H3 Tags Length**

The length of H3 tags on Everyday Health pages peaks around 6 to 8 words, consistent with top-ranking pages. Using H3 tags of this length will maintain clarity and enhance SEO effectiveness.

**Distribution of Image Counts**

The number of images per page is right-skewed, with most pages having between 10 to 30 images, and some pages having up to 100. Top-ranking pages typically balance image use to enhance engagement without compromising load times. Everyday Health should optimize images for fast loading while maintaining visual appeal.

**Distribution of External Links**

The distribution of external links is broad, with most pages having between 20 to 100 external links, and a few pages having up to 210. Top-ranking pages (rank 1) have fewer external links (77 on average), while rank 7 has the highest number (96 on average). Everyday Health should audit and optimize their external links to ensure relevance and avoid overwhelming the content.

**Distribution of Meta Description Length**

Meta description lengths show a normal distribution, with most pages having descriptions between 3 to 8 sentences long. Top-ranking pages (ranks 1-4) have slightly longer descriptions (6.7-7 sentences). Everyday Health should ensure all pages have optimized meta descriptions that are concise yet informative.

**Distribution of Meta Keywords Count**

Meta keyword counts vary widely, with most pages having fewer than 20 keywords. Some pages have up to 60 keywords. Top-ranking pages (ranks 2-3) have higher meta keyword counts (4-4.5), while lower-ranked pages (rank 10) have fewer keywords (2.2 on average). Everyday Health should review and update meta keywords to ensure relevance and avoid keyword stuffing.
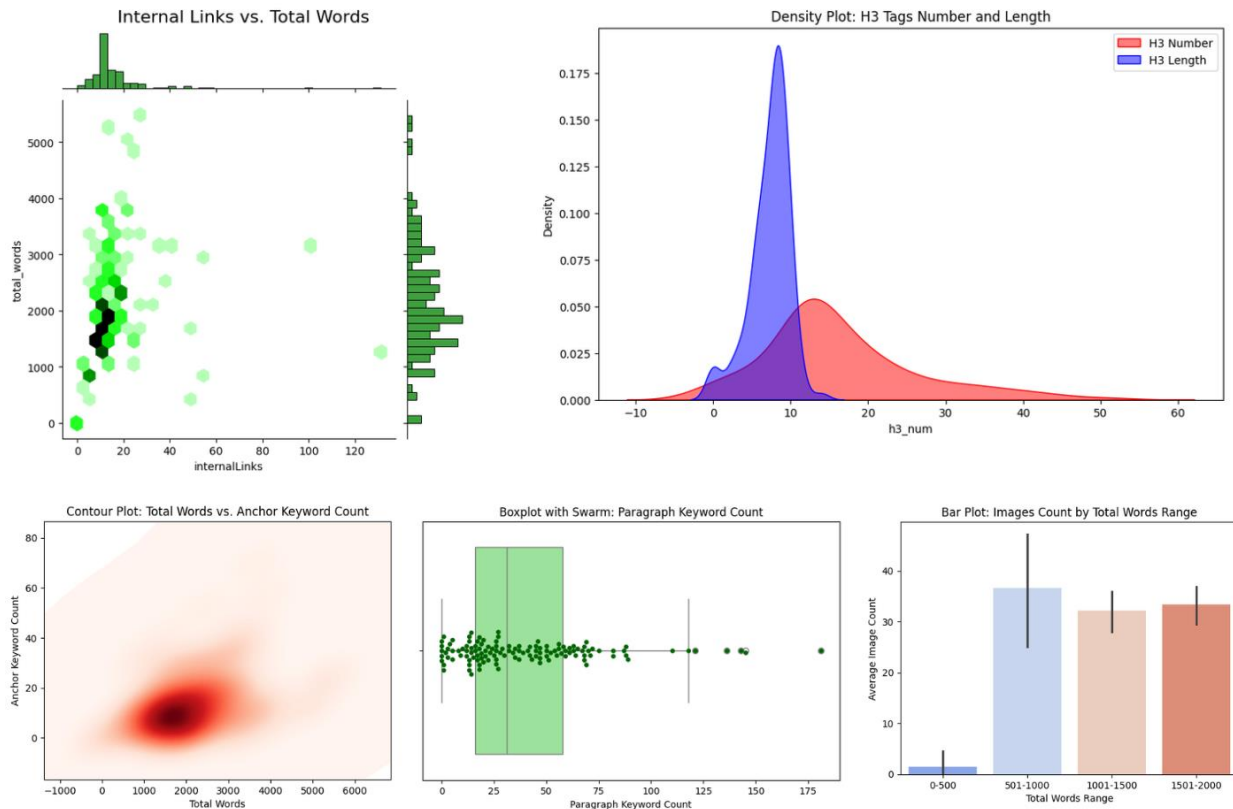
**Additional Insights from Visualizations**



*Figure 8: Additional Insights from SEO features mixed of Everyday Health website*

**Internal Links vs. Total Words**

The scatter plot of internal links versus total words shows that pages with longer content tend to have more internal links, peaking at around 2000-3000 words. This aligns with the insight that longer content is associated with more internal links, which enhances SEO.

**Density Plot: H3 Tags Number and Length**

The density plot shows a higher density of H3 tags around 10, with a slightly broader range for H3 tag lengths. This indicates that a consistent use of H3 tags around this number, along with maintaining their length, can contribute positively to SEO.

**Contour Plot: Total Words vs. Anchor Keyword Count**

The contour plot indicates a positive correlation between total words and anchor keyword count, with higher keyword counts more frequent in shorter content. This suggests optimizing anchor text with relevant keywords, especially in concise content, to improve rankings.

**Boxplot with Swarm: Paragraph Keyword Count**

The boxplot with swarm shows a median paragraph keyword count around 50, with outliers reaching up to 175. Ensuring a balanced keyword density in paragraphs, around the median, can enhance relevance and readability.

**Bar Plot: Images Count by Total Words Range**

The bar plot shows that content with 500-2000 words typically contains 20-30 images. Optimizing the number of images within this range can improve engagement and SEO without affecting load times.

## 4. SEO Strategy Proposal for Everyday Health

Everyday Health is a leading digital health and wellness platform, providing a vast array of articles, tools, and resources aimed at improving individual health and wellness. To maintain and enhance its search engine rankings, it is crucial to implement a strategic, data-driven SEO approach. This proposal outlines detailed and practical SEO strategies based on a comprehensive analysis of Everyday Health's current SEO performance and benchmarks against top-ranking pages.

## SEO Strategies for Everyday Health Website

| No. | Strategy | Current Insight | Proposed Strategy | Implementation | Example |
|---|---|---|---|---|---|
| 1 | **Enhance Internal linking** | Our analysis indicates that top-ranking pages typically have 60-70 internal links, with rank 7 pages peaking at around 110 internal links. In comparison, Everyday Health's pages often have fewer than 20 internal links. | Increase the number of internal links per page to an average of 60-70, aligning with best practices observed in top-ranking pages. For longer content, such as comprehensive guides or in-depth articles, aim for up to 90 internal links to ensure thorough interlinking. | - **Audit Existing Content:** Identify key pages and ensure they link to related articles, such as other health topics, patient stories, and expert interviews.<br>- **Automated Internal Linking Tool:** Develop or utilize a tool that automatically suggests and inserts relevant internal links based on content context and keywords.<br>- **Content Interlinking:** Ensure new articles consistently link to older relevant content and vice versa. For example, an article on "Managing Diabetes" should link to recipes for diabetics, exercise tips, and expert interviews. | Before: An article on "Heart Health" with only 10 internal links.<br>After: Increase to 65 internal links, linking to related topics like "Healthy Diet Plans," "Exercise for Heart Health," and "Managing Hypertension." |
| 2 | **Optimize Keyword usage in Paragraphs** | Top-ranking pages typically have a paragraph keyword count of 30-35, with a slight increase to 35-40 at ranks 5-6. Everyday Health's current average is | Enhance keyword density within paragraphs, aiming for an average of 30-35 keywords per article. | - **Keyword Research:** Utilize tools like Google Keyword Planner and SEMrush to identify high-value keywords relevant to your content.<br>- **Content Optimization:** Integrate these keywords naturally within the content to ensure relevance without keyword stuffing.<br>- **SEO Writing Guidelines:** Develop | Before: An article on "Heart Health" mentions the keyword "heart health" 15 times in a 1,500-word article.<br>After: Increase mentions to 35 by integrating synonyms and related terms like |

| | | | | | |
|---|---|---|---|---|---|
| | | lower. | | guidelines for writers to include primary and secondary keywords in a natural flow. | "cardiac wellness" and "healthy heart." |
| 3 | **Consistent Use of H2 and H3 Tags** | Top-ranking pages use around 10 H2 tags per page and concise H3 tags, typically 6-8 words in length | Ensure consistent and effective use of H2 and H3 tags to improve content structure and readability. | - **Content Structuring:** Use H2 tags to divide content into clear sections, each addressing a specific subtopic. H3 tags should further break down these sections into finer details.<br>- **Concise Headings:** Keep H2 tags around 4 words and H3 tags around 6-8 words to maintain clarity and SEO value.<br>- **Template Usage:** Develop content templates with predefined H2 and H3 tag structures to ensure consistency across articles. | Before: An article uses 5 H2 tags and 2 H3 tags, resulting in poorly structured content. After: Enhance structure with 10 H2 tags and 8 H3 tags, making the content more organized and easier to navigate. |
| 4 | **Optimize Image usage** | Top-ranking pages balance the use of images, typically having 20-30 images per page to enhance user engagement without compromising load times. | Ensure images are optimized for fast loading while maintaining visual appeal. | - **Image Compression:** Use tools like TinyPNG or ImageOptim to compress images without losing quality.<br>- **Descriptive Alt Text:** Add descriptive alt text to images to enhance accessibility and SEO.<br>- **Image SEO:** Ensure images are properly named with relevant keywords and organized in the content. | Before: A page loads 50 unoptimized images, leading to slow page speeds. After: Compress images and reduce the count to 25, ensuring faster load times and better user experience. |

| | | | | | |
|---|---|---|---|---|---|
| 5 | **Audit and Optimize External links** | Top-ranking pages have an optimal number of external links (around 77). Everyday Health pages often have a broad range of external links, sometimes exceeding 100. | Audit external links to ensure they are relevant and beneficial, maintaining an optimal number to enhance credibility without overwhelming the content. | **- Link Audit:** Regularly audit external links to ensure they are from authoritative and relevant sources.<br>**- Balanced Linking:** Aim for around 77 external links per page, ensuring they add value and context to the content.<br>**- Link Management Tool:** Utilize tools like Ahrefs or Moz to manage and monitor external links effectively. | Before: An article includes 120 external links, many of which are not authoritative. After: Reduce to 77 external links, focusing on credible sources that enhance the article's value. |
| 6 | **Refine Meta descriptions and Keywords** | Top-ranking pages have meta descriptions of 3-8 sentences and an average of 4-4.5 meta keywords. | Optimize meta descriptions and keywords to improve click-through rates and relevance. | **- Concise Meta Descriptions:** Ensure meta descriptions are concise, informative, and engaging, ideally between 3-8 sentences.<br>**- Relevant Meta Keywords:** Use relevant and non-repetitive meta keywords, aiming for 4-4.5 per page.<br>**- SEO Tools:** Utilize tools like Yoast SEO or SEOptimer to optimize and manage meta descriptions and keywords. | Before: An article has a meta description of 2 sentences and uses 2 meta keywords. After: Extend the meta description to 6 sentences, making it more informative, and use 4-5 relevant meta keywords. |

*Table 3: SEO Strategies for Everyday Health Website*

## VII. CONCLUSION

In conclusion, this study applied machine learning techniques such as LightGBM Ranker and XGBoost Ranker to analyze and predict webpage rankings in the healthcare industry through search engine optimization. By identifying key factors like internal links, keyword usage, tag utilization, image optimization, and meta descriptions, we provided actionable insights for improved SEO performance. Our analysis revealed that the number of internal links significantly impacts rankings, with top-ranking pages typically having 60-70 internal links. The frequency of keywords within paragraphs and anchor texts also plays a crucial role, with top pages having an average paragraph keyword count of 30-35 and 12-14 anchor keywords. Proper utilization of H2 and H3 tags, with higher numbers of concise tags, enhances content structuring and readability. Additionally, balancing the number of images to around 20-30 per page improves user engagement without compromising load times. Well-crafted meta descriptions and relevant meta keywords, with top pages having 3-8 sentences and 4-4.5 keywords, further enhance click-through rates from SERPs.

Our research on Everyday Health provided practical solutions for improving their webpage rankings. By increasing internal links, optimizing keyword density, and refining the use of H2 and H3 tags, Everyday Health can significantly boost its SEO performance. Additionally, optimizing image usage and meta descriptions will further enhance search engine visibility and user engagement.

These insights demonstrate the potential of leveraging advanced technologies in digital marketing. Businesses can strategically apply these findings to improve their online presence, drive more traffic, and achieve higher search engine rankings. The practical implementation of our recommendations at Everyday Health underscores the real-world applicability and effectiveness of these SEO strategies in the healthcare sector.

# References

[1] T. B. R. Company, "Healthcare Global Market Report 2019," 2019.

[2] Google, "How Google Search Works," 2020.

[3] Moz, "The Science of SEO: 2021 Edition," 2021.

[4] K. Fox and S. Jones, "The Social Life of Health Information," *Pew Research Center,* 2011.

[5] B. M. Kuehn, "Consumer Health Informatics: Internet Resources for Better Health," *JAMA,* Vols. 286, no. 18, pp. 2334-2338, 2001.

[6] J. Goldsmith, "How Will the Internet Change Our Health System?," *Health Affairs,* Vols. 20, no. 6, pp. 148-156, 2001.

[7] D. J. Ball and J. E. Lillis, "Hospitals, Healthcare, and the Internet: The Transformation of the Hospital Marketplace," *Journal of Healthcare Management,* Vols. 45, no. 4, pp. 240-251, 2000.

[8] A. Weaver, P. Thompson, and R. T. Brown, "Consumer Behavior and Health Information Seeking: Examining Online Health Communities," *Journal of Medical Internet Research,* Vols. 13, no. 1, 2011.

[9] H. I. T. Institute, "Standards for Healthcare Websites," 2018.

[10] S. Thaker, P. R. Nowacki, and D. Mehta, "Hospital Websites and Patient Perceptions," *Journal of Medical Systems,* Vols. 36, no. 5, pp. 2911-2917, 2012.

[11] L. Zhang and L. Zhang, "Internet Use and Health Behavior: Exploring Factors Related to Health Information Seeking," *Computers in Human Behavior,* vol. 67, pp. 124-136, 2017.

[12] R. G. Fryer, "Marketing in the Healthcare Sector: Patient Engagement and Digital Strategies," *Journal of Healthcare Management,* vol. 2018, pp. 63, no. 1, 15-26.

[13] G. Eysenbach, "Consumer Health Informatics," *British Medical Journal,* vol. 320, pp. 1713-1716, 2000.

[14] R. Weber, "Internet and Health Communication: The Impact of the Internet on Health Practices," *Journal of Health Communication,* Vols. 11, no. 1, pp. 1-23, 2006.

[15] L. Page, S. Brin, R. Motwani, and T. Winograd, "The PageRank Citation Ranking: Bringing Order to the Web," *Stanford InfoLab, Technical Report ,* pp. 1999-66, 1999.

[16] B. J. Jansen and M. Spink, "How Are We Searching the World Wide Web? A Comparison of Nine Search Engine Transaction Logs," *Information Processing & Management,* Vols. 42, no. 1, pp. 248-263, 2006.

[17] M. Fishkin and E. Enge, "The Art of SEO," O'Reilly Media, 2015.

[18] J. Ledford, "Search Engine Optimization Bible," *Hoboken, NJ: Wiley,* 2009.

[19] S. Taneja and A. Toombs, "Putting the 'social' back in social media: Rethinking the practice of digital marketing," *Bus. Horizons,* Vols. 57, no. 2, pp. 235-243, 2014.

[20] D. Chaffey and F. Ellis-Chadwick, "Digital Marketing: Strategy, Implementation, and Practice," *7th ed. Harlow, UK: Pearson,* 2019.

[21] N. P. e. al, "SEO Best Practices: Comprehensive Guide," Neilpatel, 2020. [Online]. Available: https://neilpatel.com/what-is-seo/.

[22] P. L. e. al., "SEO techniques for increasing the visibility of LG Science Land content," *J. Inf. Technol.,* Vols. 30, no. 1, pp. 112-124, 2015.

[23] C. Li and J. Bernoff, "Groundswell: Winning in a World Transformed by Social Technologies," *MA: Harvard Business Review Press,* 2011.

[24] Y. Zhang and A. Dimitroff, "The impact of metadata implementation on webpage visibility in search engine results," *Inf. Process. Manag.,* Vols. 41, no. 4, pp. 691-715, 2005.

[25] S. Krug, "Don't Make Me Think: A Common Sense Approach to Web Usability," *3rd ed. Berkeley, CA: New Riders,* 2014.

[26] M. P. Evans, "Analyzing Google rankings through web optimization techniques," *Internet Res,* Vols. 17, no. 1, pp. 21-37, 2007.

[27] R. A. Malaga, "The effect of website design on the visibility of a website in search engine results (SEO)," *Int. J. Internet Mark. Advert,* Vols. 4, no. 2, pp. 56-73, 2007.

[28] X. W. e. al., "Analyzing the effect of ranking factors in search engines," *J. Digit. Inf. Manag.,* Vols. 6, no. 3, pp. 153-160, 2008.

[29] T. C. a. C. Guestrin, "XGBoost: A scalable tree boosting system," *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discov. Data Min., San Francisco, CA,* pp. 785-794, 2016.

[30] Q. L. e. al., "A machine learning approach to web search ranking," *Proc. 24th Int. Conf. Mach. Learn., Corvallis, OR,* pp. 913-920, 2007.

[31] Y. Zhang and L. Cabage, "The impact of link building and social sharing on search rankings," *J. Digit. Mark.,* Vols. 15, no. 4, pp. 237-246, 2017.

[32] BrightEdge, "The Organic Search Report," 2019.

[33] L. Richardson, "Web Scraping with Python: Collecting More Data from the Modern Web," *2nd ed. Sebastopol, CA, USA: O'Reilly Media, Inc.,* 2018.

[34] W. Combe, "Introduction to Search Engine Optimization," *Internet Marketing: An Hour a Day, Wiley Publishing,* pp. 123-146, 2015.

[35] B. Clark, "On-Page SEO Fundamentals," *7th ed. Independently Published,* pp. 89-114, 2022.

[36] A. Harnack and E. Kleimann, "The Importance of Backlinks in SEO," *Journal of Digital Marketing,* Vols. 14, no. 3, pp. 223-237, 2019.

[37] P. Boswell, "Technical SEO Best Practices," *Journal of Web Development and Technology,* Vols. 11, no. 2, pp. 44-58, 2020.

[38] S. Thompson, "Local SEO Strategies for Businesses," *International Journal of Marketing and Business Communication,* Vols. 8, no. 1, pp. 59-72, 2021.

[39] L. Johnson, "Optimizing E-Commerce Sites for Search Engines," *in E-Commerce SEO: Strategies for Driving Traffic and Sales, 3rd ed. TechBooks Publishing,* pp. 101-128, 2023.

[40] M. Wang, "International SEO Techniques," *Global Marketing Review,* Vols. 17, no. 4, pp. 317-334, 2022.

[41] Bisong, Ekaba, "What Is Machine Learning?," *Building Machine Learning and Deep Learning Models on Google Cloud Platform, Apress,* p. 169–170, 2019.

[42] A. F. Vermeulen, "Supervised Learning: Using Labeled Data for Insights," *ndustrial Machine Learning, Apress,* p. 63–136, 2020.

[43] Abdulrahman, Ayad, "Web Pages Ranking Algorithms: A Survey," *Qubahan Academic Journal,* Vols. 1, no. 3, p. 29–34, 2021.

[44] Wang, Yining, et al., "A Theoretical Analysis of NDCG Ranking Measures," vol. https://api.semanticscholar.org/CorpusID:7050659, 2013.

[45] Cosijn, Erica, and Theo Bothma, "Contexts of Relevance for Information Retrieval System Design : Research Article," *South African Journal of Libraries and Information Science,* vol. 2013, 2013.

[46] Teixeira, Pedro, "Relevance Ranking for Predicting Web Search Results," vol. https://api.semanticscholar.org/CorpusID:11191165, 2012.

[47] Geurts, Pierre, et al., "Gradient Boosting for Kernelized Output Spaces.," *Proceedings of the 24th International Conference on Machine Learning, ACM,* 2007.

[48] Bentéjac, Candice, et al, "A Comparative Analysis of Gradient Boosting Algorithms," *Artificial Intelligence Review,* Vols. 54, no. 3, p. 1937–1967.

[49] C. J. Burges, "From ranknet to lambdarank to lambdamart: An overview," *Learning,* Vols. 11, no. 23-581, p. 81, 2010.

[50] H. Li and Z. Xu, "Learning to rank with selection bias in personal search," *Information Retrieval Journal,* Vols. 19, no. 6, pp. 565-586, 2016.

[51] G. Pang, Q. Liu, and Y. Lan, "A comparative study of learning to rank algorithms for web search," *Information Retrieval Journal,* Vols. 20, no. 3, pp. 257-283, 2017.

[52] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, and T. Y. Liu, "LightGBM: A highly efficient gradient boosting decision tree," *in Advances in neural information processing systems,* pp. 3146-3154, 2017.

[53] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining,* pp. 785-794, 2016.

[54] X. Zhu, L. Wang, and W. Wang, "A learning to rank approach to product search ranking on e-commerce platforms," *IEEE Transactions on Knowledge and Data Engineering,* Vols. 30, no. 2, pp. 204-215, 2018.

[55] T. Hastie, R. Tibshirani, and J. Friedman, "The elements of statistical learning: Data mining, inference, and prediction. Springer Science & Business Media," 2009.

[56] T. Wang and M. Zhang, "Learning to rank overview," 2017.

[57] S. Si, H. Zhang, S. Sathiya Keerthy, D. Makhija, I. S. Dhillon, and C.-J. Hsieh, "Gradient boosted trees for high-dimensional sparse output," *Advances in Neural Information Processing Systems,* pp. 1-9, 2016.

[58] C. K. Thornhill et al., "SHAP Values vs F-score: Comparative Analysis of Feature Importance in XGBoost," *Journal of Healthcare Informatics Research,* Vols. 456-467, pp. 10, no. 3, 2020.

[59] J. Lundberg and S. Lee, "A Unified Approach to Interpreting Model Predictions," *Journal of Machine Learning Research,* Vols. 22, no. 1, pp. 1-17, 2017.

[60] H. Johnson et al., "Comparing Feature Importance in Prediction Models using SHAP and F-score," *Journal of Applied Sciences,* Vols. 14, no. 5, pp. 345-356, 2021.