

# Bachelorarbeit

---

## Suchoptimierung mittels maschinellen Lernens

Zeitraum 04.07.2016 - 04.10.2016

Lukas Abegg

Matrikelnummer 798972

Sommersemester 2016

Fachsemester 6

Studiengang Medieninformatik (B.Sc.)

Beuth Hochschule für Technik



BEUTH HOCHSCHULE FÜR TECHNIK BERLIN  
University of Applied Sciences

**Beuth Hochschule für Technik**

Luxemburger Str. 10  
13353 Berlin

*Betreuer*

**Prof. Dr. habil. Alexander Löser**

Fachbereich VI - Informatik und Medien

Beuth Hochschule für Technik

*Gutachter*

**Prof. Dr. Martin Oellrich**

Fachbereich II - Mathematik - Physik - Chemie

Beuth Hochschule für Technik

# Inhaltsverzeichnis

<b>1</b>	<b>Einführung</b>	<b>1</b>
1.1	Aufbau der Suche bei Springer Nature . . . . .	1
	White Label Applikation mit Solr-Suche . . . . .	1
	User-Tracking mit Webtrekk . . . . .	1
	Architektur . . . . .	2
1.2	Problemstellung: Keine Userrelevanz in der Suche . . . . .	2
	Userrelevante Dokumente werden nicht gefunden . . . . .	2
	Der Springer Nature Stakeholder: Springermedizin setzt auf Webtrekk . . . . .	2
	Der fast gläserne User . . . . .	2
1.3	Ziel der Arbeit . . . . .	3
1.3.1	Suchoptimierung durch Click-Trough-Daten . . . . .	3
	Annahmen . . . . .	3
1.3.2	Abbildung auf das Springermedizin-Umfeld . . . . .	3
	Potential von Userrelevanzen in der Suchoptimierung analysieren . . . . .	3
	Bekanntes und wirkungsvolles Information Retrieval Verfahren . . . . .	3
	Lösungsansatz basierend auf Click-Trough-Daten aus Webtrekk . . . . .	3
	Anwendung auf Springermedizin-Umfeld . . . . .	3
	Was wird in dieser Arbeit nicht behandelt? . . . . .	4
1.4	Methodik . . . . .	4
1.4.1	Analyse der Click-Trough Daten . . . . .	4
	Userrelevante Dokumente durch Click-Trough-Rates identifizieren . . . . .	4
	Click-Trough-Daten kommen aus Webtrekk-Analysen . . . . .	4
	Suchterm semantisch aufschlüsseln mittels Segmentierung . . . . .	4
	Suchterm semantisch erweitern mittels Thesaurus . . . . .	5
	Komplexe Auswertungen der Click-Trough-Daten nicht möglich . . . . .	5
	Durch Webtrekk kein komplexer Lern-Algorithmus notwendig . . . . .	5
1.4.2	Userrelevanz in Suchprozess einbinden . . . . .	5
	Ansatz: Suchindex-Erweiterung in der Solr-Suche . . . . .	6
	Ansatz: Aufbereitung der Suchanfrage . . . . .	6
	Ansatz: Aufbereitung der Suchresultate anhand eines Klick-Modell basierten Algorithmus . . . . .	6
	Entscheidung für den Ansatz der Aufbereitung der Suchresultate anhand eines Klick-Modell basierten Algorithmus . . . . .	7
	Klick-Wahrscheinlichkeit mit Position-based Modell berechnen . . . . .	7
	Verhältnis zwischen den Klick-Wahrscheinlichkeiten abhängig der Position im Suchresultat definieren . . . . .	7
	Smoothing Faktor in Position-based Modell . . . . .	7
1.4.3	Effektive Userrelevanz . . . . .	8

	Die Klick-Wahrscheinlichkeit ist kein absoluter Wert für Userrelevanz . . . . .	8
	Overfitting vermeiden . . . . .	8
	Zusätzliche Varianz durch Zufallsfaktor . . . . .	8
1.4.4	Evaluation . . . . .	8
	Suchvarianten mithilfe eines Evaluationssystems vergleichen . . . . .	8
	Ziel der Evaluation . . . . .	8
	Evaluationssystem aufbauen . . . . .	8
	Evaluationsdaten für die Auswertung mit Cohens Kappa-Koeffizient selektieren .	9
	Evaluationsdaten mittels NDCG auswerten . . . . .	9
	Qualitätsmaß einer Suchvariante bestimmen . . . . .	9
	Verschiedene Varianten des neuen Lösungsansatzes werden evaluiert . . . . .	9
1.5	Gliederung und Aufbau . . . . .	10
	Der Lösungsansatz und deren Grundlagen . . . . .	10
	Umsetzung des Lösungsansatzes . . . . .	10
	Erkenntnisse verarbeiten . . . . .	10
<b>2</b>	<b>Grundlagen</b>	<b>11</b>
	White Label Applikation mit Solr-Suche . . . . .	11
2.1	Grundbegriffe . . . . .	11
	White Label Applikation mit Solr-Suche . . . . .	11
	<b>Literatur</b>	<b>12</b>
	<b>Abbildungs-Verzeichnis</b>	<b>13</b>
	<b>Tabellen-Verzeichnis</b>	<b>14</b>
	<b>Sourcecode-Verzeichnis</b>	<b>15</b>

# Einführung

Springer Nature ist ein weltweit führender Verlag für Forschungs-, Bildungs- und Fachliteratur mit einer breiten Palette an angesehenen und bekannten Medienmarken und zudem der weltweit größte Verlag für Wissenschaftsbücher. Für das Unternehmen Springer Nature ist es darum wichtig, auf seinen Web-Applikationen eine Suche anbieten zu können, die Suchintentionen erkennt und möglichst schnell zum gesuchten Content leitet. Die Suche wird vor allem als Hilfsmittel zur Navigation und zum Finden von Literatur und Dienstleistungen genutzt. Durch die vielen von Springer Nature publizierten Zeitschriften und Querverweise in Artikeln, wird sie aber auch oft zur Suche nach Issues<sup>1</sup> und Artikeln verwendet sowie als Hilfestellung um Diagnosen zu Krankheitsbilder stellen zu können.

Springer Nature sammelt viele User-Tracking-Daten und dadurch viel Wissen über das Verhalten der User<sup>2</sup> bei der Nutzung ihrer Suche, lässt dieses Wissen jedoch bisher noch nicht in ihre Suche einfließen. In dieser Arbeit wollen wir untersuchen, ob mithilfe dieses Wissens, die Suche optimiert werden kann.

## 1.1 Aufbau der Suche bei Springer Nature

### White Label Applikation mit Solr-Suche

Damit die verschiedenen Verlage und Zeitschriften der Verlagsgruppe Springer Nature ihre Produkte und Dienstleistungen online anbieten können nutzt Springer Nature eine eigens entwickelte White Label Applikation<sup>3</sup>. Die White Label Applikation verwendet *Apache Solr* (im Folgenden *SSolr*"genannt) (siehe [Sol]) als Suchplattform. Die Solr dient hierbei als eine der Schnittstellen zwischen dem Content-Pool von Springer Nature und der White Label Applikation. Bei den vom Content-Pool gelieferten Inhalten, handelt es sich um von Springer Nature Verlag publizierte Zeitschriften, Artikel, Bücher und redaktionelle Inhalte.

### User-Tracking mit Webtrekk

Um das Verhalten der User auf ihren Web-Applikationen zu tracken verwendet Springer Nature das Analysetool Webtrekk (siehe [Web]). Die daraus resultierenden Berichte bieten unter anderem die Möglichkeit, *Suchquery-Logs*<sup>4</sup> und *Click-Trough-Rates* (CTR)<sup>5</sup> der User auszuwerten.

---

<sup>1</sup>Nummer der Zeitschriftenausgabe, in der sich der Artikel befindet.

<sup>2</sup>Als User werden die Nutzer der Springermedizin-Suche bezeichnet

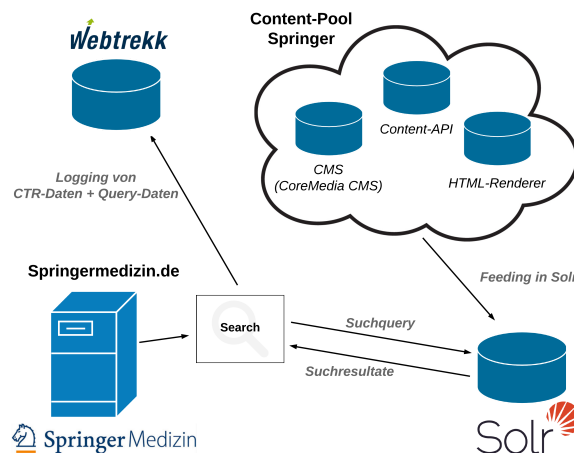
<sup>3</sup>Eine White Label Applikation ist eine wiederverwendbare und agil erweiterbare Applikation

<sup>4</sup>Protokoll über alle ausgeführten Suchanfragen auf der Applikation

<sup>5</sup>Kennzahl um die Anzahl der Klicks auf Links im Verhältnis zu den gesamten Impressionen darzustellen

## Architektur

In Abb. 1 ist die Suche nochmals grafisch aufbereitet:



**Abb. 1:** Aufbau der Suche bei Springer Nature

## 1.2 Problemstellung: Keine Userrelevanz in der Suche

### Userrelevante Dokumente werden nicht gefunden

Die User von Springermedizin suchen oft mit einschlägig, fundierten Fachbegriffen nach den neuesten und relevantesten Zeitschriften, Bücher oder Publikationen. Die zeitlich aktuellsten Suchtreffer zu finden ist für Springermedizin kein Problem. Die für den User *relevantesten* jedoch schon.

### Der Springer Nature Stakeholder: Springermedizin setzt auf Webtrekk

Zu den Stakeholder<sup>6</sup> der in Kapitel 1.1 angesprochenen White Label Applikation gehört *Springermedizin* (siehe [Sme]). Springermedizin ist ein Fortbildungs- und Informationsportal für Ärzte. Mithilfe von Webanalysten und Webtrekk versucht Springermedizin das Marketing seines Webauftrittes zu verbessern und ist sehr interessiert an neuen Ansätzen, um die gesammelten Tracking-Daten besser einzusetzen. In dieser Arbeit wird darum der Fokus auf die Verwendung von Tracking-Daten in der Suche von Springermedizin gesetzt.

### Der fast gläserne User

Springermedizin sammelt Tracking-Daten über jegliche Aktivitäten auf deren Applikationen und investiert Zeit und Geld in die Individualisierung<sup>7</sup> der Analysedaten auf Webtrekk. Mittlerweile sind knapp 30 Custom-Parameter<sup>8</sup> auf Webtrekk angelegt um genau die Daten zu tracken, die zur Analyse des Verhaltens der User auf ihrer Applikationen relevant sind. Dadurch entsteht ein fast „gläsernen User“. Dieses Wissen könnte zum Vorteil des Users eingesetzt werden, indem es in der Suche verwendet wird.

<sup>6</sup>Bezeichnet Springer Nature interne Kunden, die ein Interesse am Ergebnis der White Label Applikation haben

<sup>7</sup>Mit Individualisierung wird die Speicherung eigener Parameter bezeichnet

<sup>8</sup>Individuell erzeugte Parameter für Berichte und Analysen

## 1.3 Ziel der Arbeit

### 1.3.1 Suchoptimierung durch Click-Trough-Daten

In dieser Arbeit werden wir untersuchen, ob mithilfe der von Springermedizin gesammelten Click-Trough-Daten dessen Suche verbessert werden kann. Im Idealfall widerspiegeln die gesammelten Click-Trough-Daten der Suchresultate die Userrelevanz der einzelnen Dokumente<sup>9</sup>.

#### **Annahmen**

Wir gehen dabei von folgenden Annahmen aus. Relevante Dokumente sind wichtiger als nicht relevante Dokumente. Eine Suchergebnis ist dann gut, wenn die relevanten Ergebnisse in der verwendeten Hierarchie vor den nicht relevanten Ergebnisse auftauchen.

### 1.3.2 Abbildung auf das Springermedizin-Umfeld

#### **Potential von Userrelevanzen in der Suchoptimierung analysieren**

Die Analyse von User-Tracking-Daten bietet viel Potential bezogen auf Userrelevanzen. Sind anhand des hier umgesetzten Lösungsansatzes Verbesserungen in der Qualität der Suche zu verzeichnen, möchte Springermedizin in Zukunft vermehrt User-Tracking-Daten in die Suche einfließen lassen. Diese Arbeit könnte dann als Fundament für weitere Lösungsansätze dienen.

#### **Bekanntes und wirkungsvolles Information Retrieval Verfahren**

Suchoptimierung mittels Userrelevanz ist ein bekanntes und nicht triviales, aber relativ wirkungsvolles Information Retrieval Verfahren (siehe [EA06]). Seit Mitte der 2000er Jahre wird mithilfe dieses Verfahrens versucht, Suchmaschinen zu verbessern. Aus dieser Zeit stammen auch die ersten Ansätze um mithilfe von Click-Trough-Daten die Userrelevanz der Suchergebnisse zu berechnen (siehe [TJ05]).

#### **Lösungsansatz basierend auf Click-Trough-Daten aus Webtrekk**

Springermedizin führt ein eigenes Tracking der User durch und verwendet auf Webtrekk selbst definierte Tracking-Parameter. Dadurch hängt die Wahl des in dieser Arbeit zu untersuchenden Lösungsansatzes und dessen Umsetzung stark von den durch Webtrekk gegeben Analyse-Daten ab.

#### **Anwendung auf Springermedizin-Umfeld**

Bei der Verwendung von Userrelevanzen in der Suche handelt es sich um ein bekanntes und gut erforschtes Problem. Wir werden in dieser Arbeit versuchen, einen bestehenden Lösungsansatz (Die Wahl des Lösungsansatz folgt in Kapitel 1.4.2) auf das Springermedizin-Umfeld abzubilden. Die Herausforderung wird hierbei die Adaptierung des Lösungsansatzes auf das Springermedizin-Umfeld sein.

---

<sup>9</sup>Als Dokumente werden die einzelnen Suchresultate bezeichnet

## Was wird in dieser Arbeit nicht behandelt?

Durch den vorgegebenen Zeitraum für die Erstellung dieser Bachelorarbeit bedingt, werden wir den Lösungsansatz so wählen, dass er mit den Gegebenheiten bei Springermedizin sinnvoll und in diesem Zeitrahmen realistisch implementiert werden kann. Wir werden daher in dieser Arbeit keine Gegenüberstellung mit anderen Lösungsansätzen machen.

Bei der Umsetzung des Lösungsansatzes konzentrieren wir uns auf die Implementation des Algorithmus zur Berechnung der Userrelevanz. Die semantische Aufschlüsselung von Suchtermen ist nicht Kern dieser Arbeit. Die semantische Aufschlüsselung des Suchterms<sup>10</sup> zur Analyse der Webtrekk-Daten enthält darum keine Gewichtungen der Relationen zwischen Webtrekk-Daten und Suchterm. Alle Relationen werden gleich gewichtet.

## 1.4 Methodik

Wie in Kapitel 1.3.1 angesprochen, wollen wir das Klickverhalten der User in der Suche analysieren um die Suchergebnisse zu verbessern.

### 1.4.1 Analyse der Click-Trough Daten

#### Userrelevante Dokumente durch Click-Trough-Rates identifizieren

Eine Möglichkeit, die Userrelevanz eines Suchergebnisses zu bestimmen, ist die Verwendung der Click-Trough-Daten zur Berechnung der Userrelevanz. Anhand der Click-Trough-Rates der Suchergebnisse, können wir die für die Nutzer der Suche relevanten Dokumente identifizieren.

Dazu analysieren wir, welche Dokumente im Suchergebnis zu welchen Suchtermen angeklickt worden sind. Wir gehen dabei davon aus, dass jedes Dokument, welches zu einer Suchanfrage angeklickt wurde, eine gewisse Relevanz für diese hat.

#### Click-Trough-Daten kommen aus Webtrekk-Analysen

Als Wissensbasis zur Berechnung der Click-Trough-Daten dienen die Webtrekk-Analysen von Springermedizin. Wir können anhand dieser Analysen feststellen, welches Dokument zu welchem Suchterm angeklickt wurde. Wir wissen dadurch, wie oft es angeklickt wurde und auf welcher Position<sup>11</sup> im Suchresultat sich das Dokument dabei befunden hat.

#### Suchterm semantisch aufschlüsseln mittels Segmentierung

Ein Suchterm kann aus mehr als einem Wort bestehen. Wir müssen darum davon ausgehen, dass möglicherweise jedes Wort des Suchterms in unterschiedlichen Suchanfragen verwendet wird. Außerdem kann auch nach Synonymen eines Wortes gesucht werden. Der Suchterm muss darum semantisch aufgeschlüsselt werden, um alle relevanten Click-Trough-Daten berechnen zu können. Die Click-Trough-Daten sind dann relevant, wenn mindestens ein Wort des aufgeschlüsselten Suchterms in

<sup>10</sup>Als Suchterm wird eine Suchanfrage bezeichnet

<sup>11</sup>Mit Position wird der dargestellte Rang des Dokumentes im Suchresultat bezeichnet

Relation zu diesen Daten steht. Von einer Gewichtung der Relation wird abgesehen.

Die Auftrennung des Suchterms in die einzelne Worte können wir mithilfe einer Segmentierung<sup>12</sup> durchführen. Hier könnten wir uns überlegen, zusätzlich mit Stoppwörtern<sup>13</sup> nicht relevante Wörter aus dem Suchterm zu entfernen. Dieses Verfahren macht aber im Springermedizin-Kontext keinen sein. Wie in Kapitel 1.2 angesprochen, suchen die User der Springermedizin-Applikation oft mit einschlägig, fundierten Fachbegriffen. Wir gehen darum davon aus, dass alle Wörter des verwendeten Suchterms für das Suchergebnis relevant sind. Diese Erkenntnis basiert auf Aussagen der Redakteure von Springermedizin und Webtrekk-Analysen der meist gesuchtesten Suchtermen der letzten Monate. Auch sind Stoppwörter veraltet und werden in modernen Information Retrieval Verfahren nicht mehr eingesetzt. Wir verzichten darum auf den Einsatz von Stoppwörtern.

### Suchterm semantisch erweitern mittels Thesaurus

Für die semantische Erweiterung eines Suchwortes wird ein Thesaurus<sup>14</sup> benötigt. Die Erweiterung umfasst zum Suchterm gleichbedeutende Begriffe (Synonyme), sehr ähnliche Begriffe (Narrow Terms), ähnliche Begriffe im weiteren Sinne (Broader Terms) und verwandte Begriffe (Related Terms).

Springer Nature besitzt einen Webservice mit welchem auf den Thesaurus *Unified Medical Language System* (UMLS) (siehe [Uml]) zugegriffen werden kann.

### Komplexe Auswertungen der Click-Trough-Daten nicht möglich

Die Webtrekk-Analysen bieten uns jedoch nur beschränkte Informationen zum Klickverhalten der User. Wichtige Informationen wie die Verweildauer auf einem Dokument oder ob nach diesem Dokument ein weiteres Dokument zum gleichen Suchterm angeklickt worden ist, lassen diese Analysen nicht zu. Da diese Informationen für komplexe Auswertungen der Click-Trough-Daten nicht ausreichen, können wir in dieser Arbeit *Feedback-Strategien* für die Click-Trough-Rate Auswertung, wie in [TJ05] beschrieben, nicht verwenden. Stattdessen müssen wir davon ausgehen, dass jeder Klick auf ein Dokument relevant ist.

### Durch Webtrekk kein komplexer Lern-Algorithmus notwendig

Der Vorteil bei der Verwendung von Webtrekk ist, dass der Algorithmus nicht stetig neues Wissen lernen und altes vergessen muss. Der Algorithmus kann direkt zur Laufzeit<sup>15</sup> Analysen gegen Webtrekk über eine frei definierbare Periode machen. Dadurch kann *overfitting*<sup>16</sup> vermieden werden. Deshalb werden wir von komplexen Lern-Algorithmen wie in [EA06] vorgestellt, absehen.

## 1.4.2 Userrelevanz in Suchprozess einbinden

Es gibt drei mögliche Eingriffspunkte während des Suchprozesses, um Userrelevanzen in der Springermedizin-Suche zu verwenden. Möglich wäre die Verwendung der Userrelevanzen in der Aufbereitung der Suchanfrage auf der Springermedizin-Applikation. Denkbar wäre auch, die Berechnung der

<sup>12</sup>Bezeichnet die Aufteilung in Abschnitte, in diesem Fall in einzelne Worte

<sup>13</sup>Stoppwörter sind Wörter, die sehr häufig auftreten und für gewöhnlich keine Relevanz für den Dokumentinhalt besitzen

<sup>14</sup>Als Thesaurus wird ein strukturiertes Verzeichnis von Begriffen, welche allesamt in irgendeiner Beziehung stehen bezeichnet

<sup>15</sup>Unter Laufzeit wird in diesem Fall der Zeitpunkt der direkte Abfrage während der Suchanfrage bezeichnet

<sup>16</sup>Überanpassung des Algorithmus durch zu viele (falsche oder veraltete) Daten



Userrelevanzen in den Suchindex der Solr einzubauen. Die dritte Variante wäre die Verwendung der Userrelevanzen in der Aufbereitung der Suchresultate der Springermedizin-Applikation. Eine davon, wollen wir in dieser Arbeit untersuchen.

### **Ansatz: Suchindex-Erweiterung in der Solr-Suche**

Um die Userrelevanzen direkt in die Solr einzubeziehen gibt es zwei Varianten. Wir können das Schema des Suchindexes über die Schema API (siehe [Sch]) erweitern und alle Einträge neu indexieren, oder wir ergänzen den Index um ein externes Feld (ExternalFileField) (siehe [Ext]).

Beide Lösungsansätze ergeben nur bei der Speicherung einer einfachen Click-Count Popularität<sup>17</sup> Sinn. Diese genügen allerdings den hier gegebenen Anforderungen nicht, da die Click-Trough-Rate abhängig vom Suchterm ist. Der erste Lösungsansatz ist zudem besonders heikel, weil bei jeder Änderung des Click-Count-Wertes, das Dokument in der Solr neu indexiert werden.

### **Ansatz: Aufbereitung der Suchanfrage**

Die Solr-Suche bietet eine Boost-Funktion namens *DisMax Query Parser* (siehe [Dis]). Mit dieser können basierend auf Feldwerten, einzelne Dokumente besser im Suchergebnis positioniert werden. Die Boost-Funktion müssten wir in den Aufbau der Suchanfrage für die Suche auf der Springermedizin-Applikation einbauen. Dieser Ansatz beinhaltet einige Gefahren die wir beachten müssen.

Dazu zählen beispielsweise die Abhängigkeiten von anderen Boost-Faktoren<sup>18</sup>. Alle Boost-Faktoren hängen voneinander ab und müssten bei jeder Ergänzung um neue Faktoren normalisiert werden, um kein „über-Boosting“<sup>19</sup> einzelner Faktoren zu riskieren. Zudem besteht die Gefahr des „blinden Boosting“ von Dokumenten. Die Solr-Relevanzberechnung ist komplex und der Einfluss des „Boosting“ in die Solr-Relevanzberechnung schwer erkennbar. Auch hat Springermedizin bereits sehr schlechte Erfahrungen mit „Boosting“ gemacht und bevorzugt einen Lösungsansatz ohne „Boosting“.

### **Ansatz: Aufbereitung der Suchresultate anhand eines Klick-Modell basierten Algorithmus**

Die dritte Möglichkeit ist bei der Aufbereitung der Suchresultate aus der Solr-Suche einen Klick-Modell (siehe [Chu+15]) basierten Algorithmus zu verwenden. Dieser soll die Suchergebnisliste analysieren, die Userrelevanzen der Dokumente berechnen und die Liste neu sortieren.

Diese Lösung können wir relativ einfach in die Springermedizin-Applikation integrieren, ohne die restliche Suchlogik<sup>20</sup> zu beeinflussen. Hierbei müssen wir jedoch beachten, dass die Solr durch die Pagination-Funktion (siehe [Pag]) nur die Top-N-Ergebnisse (bei Springermedizin sind es 20 Ergebnisse) zurückgibt. Diese Logik liegt in der Springermedizin-Applikation im Aufbau der Suchanfrage. Daher können wir diese selber steuern und uns statt 20 beispielsweise die nächsten 100 Ergebnisse zurückgeben lassen. Am Ende filtern wir die ersten 20 Ergebnisse und stellen diese dar. Außerdem wissen wir bei diesem Lösungsansatz, in welcher Reihenfolge die Ergebnisse aus der Solr zurückgegeben werden. Wir kennen die Dokumente und deren Rang. Dadurch haben wir hilfreiches Zusatzwissen, welches wir in den Klick-Modell basierten Algorithmus einfließen lassen können.

<sup>17</sup>Kennzahl für alle Klicks auf ein Dokument unabhängig des Suchterms

<sup>18</sup>Die Solr besitzt eine Boosting-Funktion, um bestimmte Wertübereinstimmungen in der Suche höher gewichtet zu können

<sup>19</sup>Bezeichnet die über-priorisierte Bewertung einzelner Faktoren

<sup>20</sup>Dazu gehört die Aufbereitung der Suchanfrage für die Solr und die Suche auf der Solr

## Entscheidung für den Ansatz der Aufbereitung der Suchresultate anhand eines Klick-Modell basierten Algorithmus

Wägen wir die besprochenen Fakten ab, wirkt der Ansatz mit der Aufbereitung der Suchresultate durch einen Klick-Modell basierten Algorithmus am sinnvollsten. Wir wissen bei diesem Ansatz, welche Dokumente für die Userrelevanz-Berechnung überhaupt in Frage kommen. Zudem kennen wir alle Einfluss-Faktoren für den Algorithmus und wir sind unabhängig von der Suchlogik auf der Solr. Dadurch können wir Änderungen in unserer Logik schnell und einfach implementieren.

### Klick-Wahrscheinlichkeit mit Position-based Modell berechnen

Mithilfe der Click-Trough-Daten aus Webtrekk, können wir zwei wichtige Informationen zu jedem Suchterm ermitteln. Wir wissen welche Dokumente und welche Positionen im Suchresultat angeklickt worden sind. Zudem kennen wir die Reihenfolge der Dokumente im Suchresultat der Solr.

Ein Ansatz der genau auf diesen Informationen aufbaut, ist das *Position-based Modell* (PBM) (siehe [Chu+15]). Dieses berechnet die Wahrscheinlichkeit dafür, dass ein User ein Dokument wirklich genau analysiert, bevor er es anklickt. Es setzt sich aus zwei Wahrscheinlichkeiten zusammen. Die Wahrscheinlichkeit für einen Klick auf die Position im Suchresultat und die Wahrscheinlichkeit für einen Klick auf das Dokument. Diesen Ansatz werden wir in dieser Arbeit implementieren.

### Verhältnis zwischen den Klick-Wahrscheinlichkeiten abhängig der Position im Suchresultat definieren

Aus eigener Erfahrung wissen wir, dass die ersten Dokumente im Suchresultat immer zuerst gesehen werden. Die dahinter gelisteten Dokumente werden fortlaufend analysiert. Wir sollten darauf achten, dass je *schlechter* der Rang des angeklickten Dokumentes im Suchresultat der Solr ist, desto *höher* das Relevanzfeedback zu bewerten ist.

### Smoothing Faktor in Position-based Modell

Wir wissen dass eine Wahrscheinlichkeit einen Wert zwischen 1 und 0 besitzt. Dadurch können Nullwerte entstehen. Das PBM multipliziert die Positions- und Dokument-Wahrscheinlichkeit miteinander, um die Klick-Wahrscheinlichkeit zu berechnen. Wir müssen aber davon ausgehen, dass es Dokumente geben kann, deren Rang nie angeklickt worden ist und umgekehrt.

Multiplikationen mit Null ergeben immer einen Nullwert. An dieser Stelle führen wir einen *Smoothing-Faktor* ein. Der Smoothing-Faktor soll zwei Probleme lösen. Zum einen wollen wir einen Wahrscheinlichkeitswert trotz der Multiplikation mit Null beachten. Zum anderen wollen wir die im vorherigen Absatz beschriebene Gewichtung abhängig des Relevanzfeedbacks in den Algorithmus einbeziehen. Wir transformieren dazu das Produkt der beiden Wahrscheinlichkeiten in eine gewichtete Summe, dem sogenannten *Weighted Moving Average*, dessen Gewichte sich zu Eins aufsummieren. Diese Gewichte sind die Smoothing-Faktoren, weshalb das Verfahren zählt zu den Smoothing-Algorithmen zählt.

### 1.4.3 Effektive Userrelevanz

#### Die Klick-Wahrscheinlichkeit ist kein absoluter Wert für Userrelevanz

Nun könnten wir die Klick-Wahrscheinlichkeit als absoluten Wert für die *Userrelevanz* betrachten. Dies wäre jedoch falsch, wir müssen davon ausgehen, dass viele User der Qualität der Suchmaschine vertrauen. Diese betrachten die Top-Suchresultate als die relevanten Suchresultate. Denkbar wäre auch, dass sie unabsichtlich das falsche Dokument anklicken.

#### Overfitting vermeiden

Um dem entgegenzuwirken und ein *overfitting* zu vermeiden, darf der Algorithmus nicht immer anschlagen. Wir müssen sicherstellen, dass vereinzelt zufällige Dokumente in den „Top-Suchresultaten“ angezeigt werden. So können auch andere Dokumente in den Fokus des Users gerückt werden. Das System fährt sich dadurch nicht auf falschen Annotationen fest.

#### Zusätzliche Varianz durch Zufallsfaktor

Mithilfe eines Zufallsfaktors kann eine solche Varianz in den Klick-Modell basierten Algorithmus gebracht werden. Wie bereits weiter oben erwähnt, werden viele Suchresultate nie und deren Rang selten bis gar nicht angeklickt. Sie haben darum keine Click-Trough-Daten. Deren Klick-Wahrscheinlichkeit ist entweder Null oder sehr klein. Der Zufallsfaktor soll darum nur leichte Einflüsse in die Klick-Wahrscheinlichkeitsberechnung haben. Auch hier können wir wieder mit dem oben eingeführten *Weighted Moving Average* (siehe [Kar]) arbeiten.

### 1.4.4 Evaluation

#### Suchvarianten mithilfe eines Evaluationssystems vergleichen

Das große Kernproblem der Überprüfung der Verbesserungen durch den untersuchten Lösungsansatz wird das Messen der Qualität der erzielten Suchergebnisse sein. Mithilfe einer Evaluation wollen wir messen, wie gut die Suchergebnis-Qualität der aktuellen Springermedizin-Suche im Vergleich zur im Zuge dieser Arbeit entwickelten Lösung ist.

#### Ziel der Evaluation

Die Evaluation soll Informationen darüber liefern, wie viel Verbesserung der neue Lösungsansatz bringt. Aus den Ergebnissen wollen wir erkennen, an welchen „Schrauben“ etwas gedreht werden muss, damit die Suche wirklich gute Ergebnisse aus Sicht der User bringt.

#### Evaluationssystem aufbauen

Dazu müssen wir eine passende Testumgebung aufbauen. Diese besteht aus einem Evaluationssystem, einer Instanz der aktuellen Springermedizin-Applikation und einer Instanz des neu implementierten Lösungsansatzes. Auf dem Evaluationssystem sollen fachliche Experten (Redakteure von Springermedizin) die Relevanz der Suchergebnisse der beiden Suchmaschinen vergleichen. Dazu sollen die jeweils

besten 10 Suchergebnisse nach Relevanz zum Suchterm bewertet werden. Der Ergebnisse werden in einer Datenbank gespeichert, um sie später auszuwerten.

### **Evaluationsdaten für die Auswertung mit Cohens Kappa-Koeffizient selektieren**

Um die Zuverlässigkeit der Relevanzbewertungen zu messen, werden wir die gleichen Suchterme jeweils von zwei fachlichen Experten bewerten lassen. Das meist verwendete Maß zur Bewertung der Übereinstimmungsgüte ist der *Cohens Kappa-Koeffizient* (siehe [Gro+07]). Diese Zahl misst den Anteil übereinstimmender Bewertungen. Hierbei können aber auch zufällige Übereinstimmungen entstehen. Der Cohens Kappa-Koeffizient korrigiert das Maß an Übereinstimmung um diesen Zufallsfaktor. Anhand der Auswertungen werden wir ein Mindestmaß der Übereinstimmungsgüte definieren. Die darunter liegenden Bewertung werden wir in der Auswertung ignorieren.

### **Evaluationsdaten mittels NDCG auswerten**

Um das Qualitätsmaß der beiden Suchen vergleichen zu können werden wir den Bewertungsalgorithmus *NDCG* (siehe [Wan+13]) einsetzen. Dieser geht davon aus, dass besser positionierte Suchergebnisse eine höhere Relevanz als schlechter positionierte haben. Der NDCG vergleicht die Reihenfolge der Relevanzbewertungen der Suchergebnisse mit der idealen Reihenfolge derselben Relevanzbewertungen. Im Idealfall entspricht die Reihenfolge der Suchergebnisse der Relevanz der Suchergebnisse.

### **Qualitätsmaß einer Suchvariante bestimmen**

In der Evaluation werden zu jedem Suchterm zwei Bewertungen für die Springermedizin-Suche und zwei Bewertungen für die Suche mit dem hier zu untersuchenden Lösungsansatz abgegeben. Um das Qualitätsmaß einer Suchvariante zu einem Suchterm zu bestimmen, berechnen wir den NDCG der beiden Bewertungen. Nehmen wir den Mittelwert der beiden resultierenden NDCG-Werte, erhalten wir den effektiven NDCG-Wert. Die NDCG-Werte der beiden Suchen können wir dann miteinander vergleichen.

### **Verschiedene Varianten des neuen Lösungsansatzes werden evaluiert**

Der in dieser Arbeit zu untersuchende Lösungsansatz kann verschieden konfiguriert werden. Wir können den Einfluss des Zufallsfaktors bestimmen. Um verschiedene Konstellationen testen zu können, werden wir mit zwei verschiedenen Werten für den Einfluss des Zufallsfaktor evaluieren. Die Click-Trough-Daten von an der Applikation angemeldeten Benutzern können wir von den Click-Trough-Daten von anonymen Benutzern unterscheiden.

Aus den beiden Einflusswerten des Zufallsfaktors und der Unterscheidung zwischen angemeldeten und anonymen Benutzern, ergeben sich vier Konstellationen, die evaluiert werden können. Jeder Konstellation werden wir jeweils 25 Prozent der Suchterme zuteilen. Mithilfe des Evaluationssystems werden wir die Zuteilung der Suchterme zufällig generieren lassen.

## 1.5 Gliederung und Aufbau

### **Der Lösungsansatz und deren Grundlagen**

Im ersten Kapitel wurde der zu untersuchende Lösungsansatz vorgestellt. Dabei sind wir auf die Hintergründe dieser Arbeit und die Vorgehensweise eingegangen. Im zweiten Kapitel (Grundlagen) folgt die Theorie des beschriebenen Lösungsansatzes. Hier werden wir uns auf die fachlichen Grundlagen konzentrieren.

### **Umsetzung des Lösungsansatzes**

In Kapitel 3 (Reranking mittels Click-Trough-Rate Ergebnis) werden wir die in Kapitel 1.4 angesprochene Methodik verfeinern und detailliert die Vorgehensweise bei der Umsetzung diskutieren. Die Umsetzung selbst folgt dann in Kapitel 4 (Implementierung).

### **Erkenntnisse verarbeiten**

Um zu prüfen ob der umgesetzte Lösungsansatz die erhofften Verbesserungen erzielt, werden wir diesen in Kapitel 5 (Evaluation und Auswertung) in einer Evaluation mit der bisherigen Springermedizin-Suche vergleichen. Aufgrund der resultierenden Erkenntnisse werden wir in Kapitel 6 ein Fazit ziehen können und einen Ausblick auf mögliche zukünftige Arbeiten geben.

# Grundlagen

White Label Applikation mit Solr-Suche

## 2.1 Grundbegriffe

White Label Applikation mit Solr-Suche

# Literatur

- [Chu+15] Aleksandr Chuklin, Ilya Markov und Maarten de Rijke. *Click Models for Web Search*. Morgan & Claypool, 2015 (siehe S. 6, 7).
- [Dis] *The DisMax Query Parser - Apache Solr Reference Guide* (siehe S. 6).
- [EA06] Susan Dumais Eugene Agichtein Eric Brill. *Improving Web Search Ranking by Incorporating User Behavior Information*. Report. Microsoft, 2006 (siehe S. 3, 5).
- [Ext] *ExternalFileField (Solr 4.8.0 API) - Apache Lucene* (siehe S. 6).
- [Gro+07] Grouven, Bender, Ziegler und Lange. „Der Kappa-Koeffizient“. In: Thieme, 2007, S. 111–128 (siehe S. 9).
- [Kar] Marzena Narodzonek Karpowska. *Smoothing methods* (siehe S. 8).
- [Pag] *Pagination of Results - Apache Solr Reference Guide* (siehe S. 6).
- [Sch] *Schema API - Apache Solr Reference Guide - Apache Software* (siehe S. 6).
- [Sme] *Springer Medizin – das Fachportal für Ärzte* (siehe S. 2).
- [Sol] *Solr - Apache Lucene* (siehe S. 1).
- [TJ05] Bing Pan Thorsten Joachims Laura Granka. *Accurately Interpreting Clickthrough Data as Implicit Feedback*. Report. Random, 2005 (siehe S. 3, 5).
- [Uml] *Unified Medical Language System (UMLS)* (siehe S. 5).
- [Wan+13] Yining Wang, Liwei Wang, Yuanzhi Li u. a. „A Theoretical Analysis of NDCG Type Ranking Measures“. In: CoRR abs/1304.6480 (2013) (siehe S. 9).
- [Web] *Webtrekk* (siehe S. 1).

# Abbildungs-Verzeichnis

1	Aufbau der Suche bei Springer Nature . . . . .	2
---	--	---



## Tabellen-Verzeichnis

## Sourcecode-Verzeichnis