

# Bachelorarbeit

---

## Suchoptimierung mittels maschinellen Lernens

Zeitraum 04.07.2016 - 11.10.2016

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
	Suchterm semantisch aufschlüsseln . . . . .	4
	Aufbereitung Click-Trough-Daten . . . . .	4
	Click-Trough-Rate in Suchprozess einbinden . . . . .	4
	Result-Reranking mittels PBM Algorithmus . . . . .	5
	Vergessen der alten Daten . . . . .	5
1.5	Gliederung und Aufbau . . . . .	5
	Der Lösungsansatz und deren Grundlagen . . . . .	5
	Umsetzung des Lösungsansatzes . . . . .	5
	Erkenntnisse verarbeiten . . . . .	5
<b>2</b>	<b>Grundlagen</b>	<b>6</b>
2.1	Grundbegriffe . . . . .	6
2.1.1	Semantik von User-Interaktionen . . . . .	6
	Synonyme und verwandte Begriffe . . . . .	6
	Wenige Dokumente erhalten viele Klicks . . . . .	6
	Niedrige Positionen werden häufiger angeklickt . . . . .	7
2.1.2	Userrelevanz mittels Click-Trough-Rate (CTR) . . . . .	8
	Was sind Click-Trough-Daten und wie entstehen diese? . . . . .	8
	Wie werden die Click-Trough-Daten in Webtrekk gespeichert? . . . . .	8

Wie können wir Click-Trough-Daten aus Webtrekk lesen? . . . . .	9
Wie sehen die Click-Trough-Daten aus? . . . . .	9
Das Userverhalten bestimmt über die Relevanz eines Klicks . . . . .	9
2.1.3 Result-Reranking mittels PBM Algorithmus . . . . .	9
2.2 Zusammenfassung . . . . .	10
<b>3 Reranking mittels Click-Trough-Rate Ergebnis</b>	<b>11</b>
3.1 Probleme und Lösungsansätze . . . . .	11
3.1.1 Suchterm Segmentierung . . . . .	11
3.1.2 Aufbereitung Click-Trough-Daten . . . . .	12
Kein Einfluss auf Suchergebnisqualität während der Klicks . . . . .	12
Userverhalten . . . . .	12
3.1.3 Click-Trough-Rate in Suchprozess einbinden . . . . .	12
3.1.4 Result-Reranking mittels PBM Algorithmus . . . . .	12
3.2 Methodik . . . . .	12
3.2.1 Suchterm Segmentierung . . . . .	12
Suchterm semantisch aufschlüsseln mittels Segmentierung . . . . .	12
Suchterm semantisch erweitern mittels Thesaurus . . . . .	13
3.2.2 Aufbereitung Click-Trough-Daten . . . . .	13
Jeder Klick auf ein Dokument ist relevant . . . . .	13
Gewichtung der Click-Trough-Daten . . . . .	13
Berechnung der Click-Trough-Rate . . . . .	13
3.2.3 Click-Trough-Rate in Suchprozess einbinden . . . . .	14
Ansatz: Suchindex-Erweiterung in der Solr-Suche . . . . .	14
Ansatz: Aufbereitung der Suchanfrage . . . . .	14
Ansatz: Aufbereitung der Suchresultate anhand eines Klick-Modell basierten Algorithmus . . . . .	14
Entscheidung für den Ansatz der Aufbereitung der Suchresultate anhand eines Klick-Modell basierten Algorithmus . . . . .	15
3.2.4 Result-Reranking mittels PBM Algorithmus . . . . .	15
Klick-Wahrscheinlichkeit mit Position-based Modell berechnen . . . . .	15
Verhältnis zwischen den Klick-Wahrscheinlichkeiten abhängig der Position im Suchresultat definieren . . . . .	15
Smoothing Faktor in Position-based Modell . . . . .	16
3.2.5 Vergessen der alten Daten . . . . .	16
Durch Webtrekk ist kein komplexer Lern-Algorithmus notwendig . . . . .	16
Die Klick-Wahrscheinlichkeit ist kein absoluter Wert für die Userrelevanz . . . . .	16
Overfitting vermeiden . . . . .	17
Zusätzliche Varianz durch Zufallsfaktor . . . . .	17
3.3 Der PBM-Algorithmus . . . . .	17
3.4 Zusammenfassung . . . . .	17
<b>4 Implementierung</b>	<b>18</b>
4.1 Technologie-Stack . . . . .	18
4.2 Architektur der Implementierung . . . . .	18
4.3 Highlight: Webtrekk-Analysen . . . . .	18
4.4 Highlight: PBM Rerank-Algorithmus . . . . .	18
4.5 Zusammenfassung . . . . .	18

<b>5</b>	<b>Evaluation und Auswertung</b>	<b>19</b>
5.1	Einführung . . . . .	19
	Suchvarianten mithilfe eines Evaluationssystems vergleichen . . . . .	19
	Ziel der Evaluation . . . . .	19
5.2	Aufbau der Analyse . . . . .	19
5.2.1	Datengrundlage . . . . .	19
	Filterung der nutzbaren Daten mittels Cohens Kappa . . . . .	19
5.2.2	Metrik . . . . .	19
	Evaluationsdaten mittels NDCG-Algorithmus auswerten . . . . .	19
	Qualitätsmaß einer Suchvariante bestimmen . . . . .	20
5.2.3	Vorgehen . . . . .	20
	Evaluationssystem aufbauen . . . . .	20
	Evaluationssystem auswerten . . . . .	20
5.2.4	Durchführung . . . . .	20
	Verschiedene Varianten des neuen Lösungsansatzes werden evaluiert . . . . .	20
5.3	Auswertung der Suchergebnis-Qualität . . . . .	21
5.3.1	Quantitative Auswertung . . . . .	21
5.3.2	Diskussion . . . . .	21
5.4	Zusammenfassung . . . . .	21
<b>6</b>	<b>Zusammenfassung und Ausblick</b>	<b>22</b>
6.1	Zusammenfassung . . . . .	22
6.2	Ausblick . . . . .	22
<b>7</b>	<b>Anhang</b>	<b>23</b>
	Literatur-Verzeichnis . . . . .	23
	Abbildungs-Verzeichnis . . . . .	24
	Tabellen-Verzeichnis . . . . .	25

# 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 [solr]) 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 [webtrekk]). 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 Springer Nature Suche bezeichnet

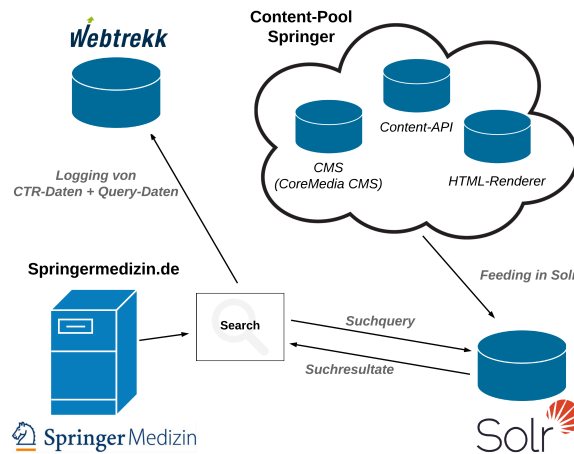
<sup>3</sup>Eine White Label Applikation ist eine wiederverwendbare und agil erweiterbare Applikation

<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 Abbildung 1 ist die Suche nochmals grafisch aufbereitet:



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

## 1.2 Problemstellung: Keine Userrelevanz in der Suche

Zu den Stakeholder<sup>6</sup> der in Kapitel 1.1 angesprochenen White Label Applikation gehört *Springermedizin* (siehe [SMED]). Springermedizin betreibt ein Fortbildungs- und Informationsportal für Ärzte.

### 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

Mithilfe von Webanalysten und Webtrekk versucht Springermedizin das Marketing seines Webauftritts 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. Konkret wollen wir dies anhand des Klick-Modell basierten Algorithmus *Position-Based Modell* (siehe [pbm]) untersuchen. 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. Ein Suchergebnis ist dann gut, wenn die relevanten Ergebnisse in der verwendeten Hierarchie vor den nicht relevanten Ergebnissen 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 [IWUSBI]). 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 [Joachims]).

#### 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 gegebenen Analyse-Daten ab.

#### Anwendung auf Springermedizin-Umfeld

Wir werden in dieser Arbeit versuchen, Click-Trough-Rates mithilfe eines Algorithmus basierend auf dem Klick-Modell *Position-Based Modell* (siehe [pbm]) zu berechnen und mit diesen ein *Reranking*<sup>10</sup> der Suchresultate 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

<sup>10</sup>Mit Reranking bezeichnen wir die Umsortierung einer Liste von Suchresultaten

## 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.

## 1.4 Methodik

Wie in Kapitel 1.3.1 angesprochen, wollen wir das Klick-Verhalten der User in der Suche analysieren, um mithilfe der daraus berechenbaren Click-Trough-Rates, die Suchergebnisse zu verbessern. Dieses Klick-Verhalten können wir aus den Click-Trough-Daten lesen.

### Suchterm semantisch aufschlüsseln

Um mit den Click-Trough-Daten arbeiten zu können, müssen wir zunächst die relevanten Click-Trough-Daten herausfiltern. Dazu müssen wir die Click-Trough-Daten dem *Suchterm* Als Suchterm wird eine Suchanfrage bezeichnet der Anfrage zuordnen können. Zu den Click-Trough-Daten wird immer der Suchterm gespeichert, mit dem dabei gesucht wurde. Das heißt wir können eine Relation zwischen dem Suchterm der Click-Trough-Daten und dem Suchterm unserer Anfrage herstellen.

Die Click-Trough-Daten müssen aber nicht mit dem vollständigen Suchterm in Relation stehen. Sie können auch nur mit einem Wort, einem Teil des Suchterms oder einem Synonym eines dieser Worte in Relation stehen. Wir müssen darum den Suchterm semantisch aufschlüsseln um alle relevanten Click-Trough-Daten filtern zu können.

### Aufbereitung Click-Trough-Daten

Können wir alle relevanten Click-Trough-Daten zu einer Suchanfrage filtern, müssen lernen wie wir diese richtig aufbereiten, um die Click-Trough-Daten berechnen zu können. Ein wichtiger Punkt bei der Aufbereitung der Click-Trough-Daten ist die Interpretation des Relevanzfeedbacks der einzelnen Click-Trough-Daten. Nicht jeder Klick ist gleich relevant zu interpretieren. Die Relevanz eines Klicks hängt davon ab, welche Aktionen der User während dem Suchvorgang vor und nach dem Klick durchgeführt hat. Wir müssen darum zuerst analysieren, welche Informationen wir zu den Click-Trough-Daten aus Webtrekk lesen können. Reichen die Informationen für detailliertere Interpretationen nicht aus, müssen wir alle Click-Trough-Daten als gleich relevant lesen.

### Click-Trough-Rate in Suchprozess einbinden

Wissen wir, wie wir die Click-Trough-Daten aufbereiten, können wir diese zur Berechnung der Click-Trough-Rate eines Dokumentes und somit dessen Userrelevanz einsetzen. Für die Verwendung dieser Userrelevanz müssen wir festlegen wann und in welcher Art wir sie einsetzen. Hierfür werden wir verschiedene Eingriffspunkte während des Suchprozesses analysieren.

Wie bereits in Kapitel 1.3.2 erwähnt, werden wir uns in dieser Arbeit auf die *Aufbereitung* der Suchresultate aus der Solr konzentrieren und dort einen Reranking-Algorithmus einbauen. Dazu



werden wir zuerst anschauen, warum wir mit diesem Ansatz arbeiten und was gegen die alternativen Eingriffspunkte spricht.

### **Result-Reranking mittels PBM Algorithmus**

Der im vorherigen Abschnitt erwähnte Algorithmus basiert auf einem Klick-Modell. Er verwendet die Click-Trough-Daten, um daraus die Click-Trough-Rate eines Dokumentes zu berechnen. Die Wahl des Klick-Modells hängt darum von den Click-Trough-Daten und den darin enthaltenen Informationen ab. In dieser Arbeit werden wir veranschaulichen, warum wir den Ansatz des Position-Based Modells (siehe [pbm]) gewählt haben und wie wir den Algorithmus umgesetzt haben.

### **Vergessen der alten Daten**

Das Position-Based Modell berechnet Wahrscheinlichkeiten, um die Click-Trough-Rate eines Dokumentes zu einer Suchanfrage festzustellen. Dem Algorithmus muss dazu das notwendige Wissen entweder zur Verfügung gestellt oder antrainiert werden. Wird dem Algorithmus das Wissen antrainiert, benötigt der Algorithmus eine Möglichkeit, neues Wissen zu Lernen und altes zu Vergessen.

Mit Webtrekk haben wir eine Wissensbasis, die durch die Springermedizin-Applikation automatisch um neue Click-Trough-Daten ergänzt wird. Das heißt wir müssen uns nicht um eine Möglichkeit zum Lernen neuer Daten kümmern. Wir müssen uns aber überlegen, wie wir altes Wissen vergessen und wie wir dem User neues Wissen präsentieren, damit dieser sich durch die Click-Trough-Rate-Berechnung nicht auf alten Dokumenten festfährt.

## **1.5 Gliederung und Aufbau**

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

In diesem Kapitel wurde der zu untersuchenden 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 (Zusammenfassung und Ausblick) ein Fazit ziehen können und einen Ausblick auf mögliche zukünftige Arbeiten geben.

# Grundlagen

## 2.1 Grundbegriffe

In diesem Kapitel werden wir die fachlichen Grundlagen zu unserem in Kapitel 1.4 vorgestellten Lösungsansatz aufarbeiten. Wichtig hierfür ist das Verständnis für die Problemstellungen in der Interaktion zwischen den Nutzern der Suche und der Suche selbst. Dazu gehört die Auseinandersetzung mit dem Klick-Verhalten der User auf der Suche von Springermedizin. Mit der Click-Trough-Rate wollen wir eine Userrelevanz bestimmen. Dazu müssen wir die Click-Trough-Daten als Relevanzfeedback deuten können. Wie wir die Click-Trough-Rates berechnen, lernen wir in den Grundlagen zu unserem Reranking-Algorithmus. Diese Grundlagen sind notwendig für die Umsetzung des Lösungsansatzes in den nachfolgenden Kapiteln 3 und 4.

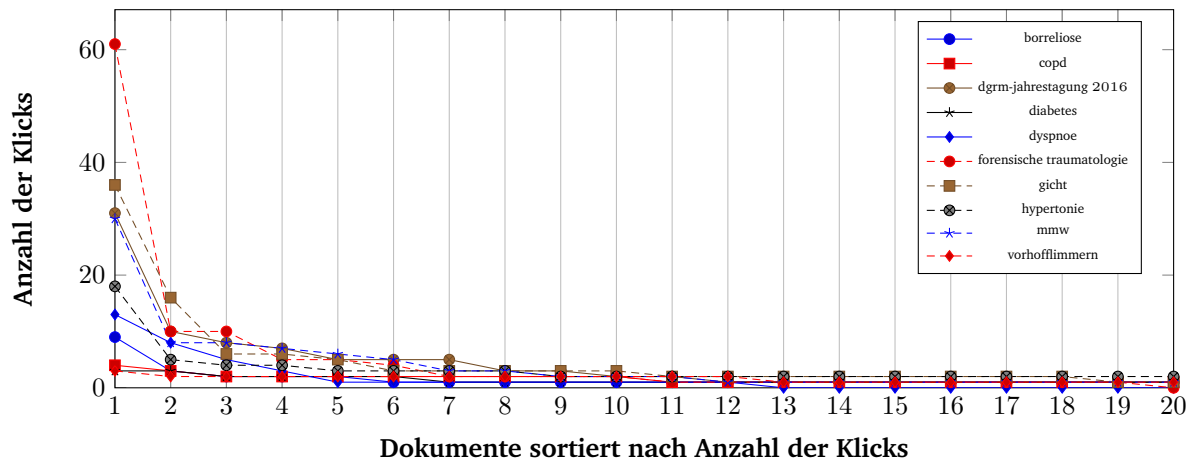
### 2.1.1 Semantik von User-Interaktionen

#### Synonyme und verwandte Begriffe

Nehmen wir als Beispiel die Suchanfrage „chronische Dyspnoe“. Würde eine andere Person stattdessen nach dem gleichbedeutenden Suchterm „konstante Atemnot“ suchen, könnten auf derer Bedeutungs-gleichheit die Ergebnisse beider Suchresultate und somit auch deren Click-Trough-Daten gleichermaßen relevant sein.

#### Wenige Dokumente erhalten viele Klicks

Um das Klick-Verhalten der User auf der Springermedizin-Suche zu verstehen, ist es wichtig anhand oft gesuchter Suchphrasen dieses Verhalten zu analysieren. Dazu wurde eine Analyse über einen Zeitraum von 30 Tagen erstellt und die zehn am häufigsten gesuchten Suchphrasen verwendet. Die Analyse vergleicht für jede Suchphrase die 20 Dokumente mit den meisten Klicks. Die Dokumente wurden hierbei nicht nach Position im Suchergebnis sondern nach Klick-Wert selektiert. Jeder Graph stellt eine Suchphrase dar. Die meisten Graphen zeigen ein stark exponentiell abnehmendes Verhalten. Dieses exponentielle Verhalten zeigt, dass einzelne Dokumente häufig und viele Dokumente selten bis nie angeklickt werden. Dieser Effekt kann wie in vielen natürlichen Phänomenen, durch das Potenzgesetz (Power Law, siehe [PowerLaw]) beschrieben werden.



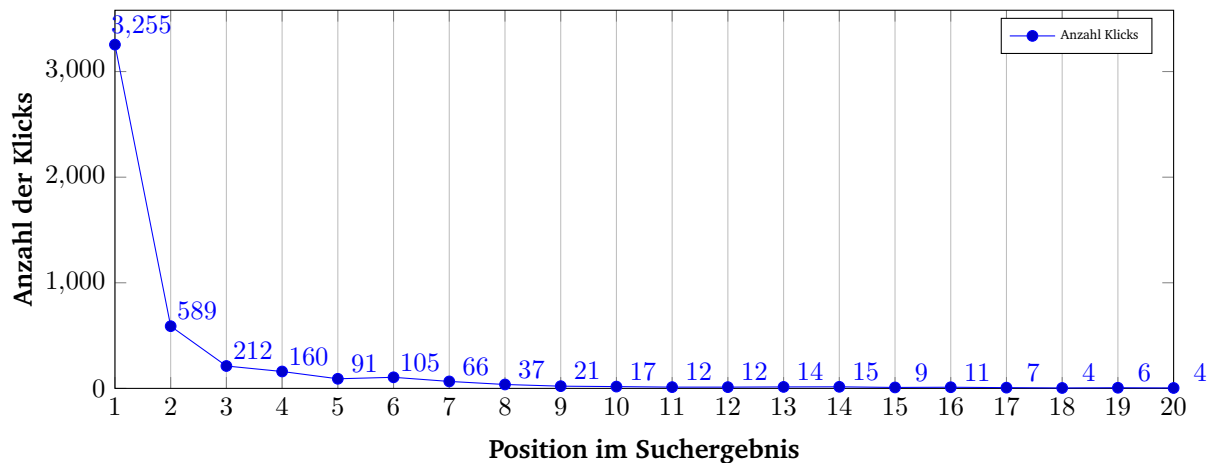
**Abb. 2:** Analyse der 20 am häufigsten angeklickten Dokumente der zehn meistgesuchten Suchphrasen.  
Zeitraum der Analyse: 19.08.16 - 19.09.16

Betrachten wir die Graphen können wir vor allem auf der ersten Position hohe Klick-Werte feststellen. Daraus lässt sich vermuten, dass einzelne Dokumente eine sehr hohe Relevanz für die entsprechende Suchanfrage aufweisen und nur wenige Dokumente auf die User als relevant wirken. Ein weitere Vermutung könnte sein, dass der zu durchsuchende Content wenig relevante Dokumente hat. Die Suchphrasen lassen auf sehr diverse Suchintentionen deuten. Es handelt sich hierbei unter anderem um Krankheiten, Zeitschriften und Behandlungen mit mehreren tausend Suchergebnissen. Die Wahrscheinlichkeit, dass wenig relevanter Content für die meisten der analysierten Suchphrasen zutrifft, sollte deshalb relativ gering sein. Wir müssen darum eher davon ausgehen, dass sich die User auf einzelne Dokumente festfahren. Das könnte an schlechten Suchergebnissen und somit an einer schlechten Suchqualität liegen. Um jedoch ein genaueres Bild über das Verhalten erstellen zu können müssen wir einen Vergleich mit der nachfolgenden Analyse in Abbildung 3 ziehen.

### Niedrige Positionen werden häufiger angeklickt

In der unten folgenden Analyse sehen wir das positionsbezogene Klick-Verhalten der User der Springermedizin-Suche. Dazu wurden über den Zeitraum von einem Monat, die letzten 1000 Suchanfragen ausgewertet. Dargestellt sehen wir die Häufigkeitsverteilung der Klicks als Graph. Wir beschränken uns hierbei auf die ersten 20 Positionen der Suchresultate. Wie wir sehen, nimmt die Anzahl der Klicks mit zunehmender Position exponentiell ab. Dieser Effekt kann ebenfalls wie in Abbildung 2, durch das Potenzgesetz (Power Law, siehe [PowerLaw]) beschrieben werden.

Betrachten wir den Graph, sehen wir, dass besonders die erste Position auffällig viel angeklickt wird. Daraus könnten wir die Vermutungen ziehen, dass die Suche eine sehr gute Qualität besitzt und die meisten User der Suchmaschine vertrauen. Wie wir aber aus den Analysen von [Joachims] lesen können, müssen wir eher davon ausgehen, dass die meisten User der Suchmaschine vertrauen und darum größtenteils ein Dokument aus den ersten zehn Positionen anklicken. Vergleichen wir die Analyse aus Abbildung 2 mit dieser Analyse, sehen wir ein sehr ähnliches Muster in der Häufigkeitsverteilung der Klicks. Wir können anhand der Klick-Zahlen ebenfalls vermuten, dass die am häufigsten angeklickten Dokumente sich dabei auf den kleinsten Positionen befunden haben.



**Abb. 3:** Analyse der Klicks auf die ersten 20 Positionen der Suchergebnisse aller Suchanfragen.  
Zeitraum der Analyse: 19.08.16 - 19.09.16

## 2.1.2 Userrelevanz mittels Click-Trough-Rate (CTR)

Um mit Click-Trough-Daten arbeiten zu können, müssen wir zuerst verstehen, was Click-Trough-Daten sind und wie sie entstehen.

### Was sind Click-Trough-Daten und wie entstehen diese?

Click-Trough-Daten sind Tracking-Daten. Tracking-Daten entstehen durch die Interaktion zwischen dem User der Applikation und der Applikation selbst. Sie verfolgen das Verhalten der User auf der Applikation und speichern diese in einer Datenbank, in unserem Fall in Webtrekk ab. Die für uns interessanten Tracking-Daten entstehen, wenn der User auf der Suche von Springermedizin ein Anfrage stellt und darauf folgend, ein Element aus dem Suchresultat anklickt.

### Wie werden die Click-Trough-Daten in Webtrekk gespeichert?

Die Speicherung der Daten auf Webtrekk übernimmt die Springermedizin-Applikation. Führt ein User eine Suche durch und klickt dabei ein Resultat an, sendet die Springermedizin-Applikation die Tracking-Informationen an Webtrekk. Die Tracking-Daten für diese Aktion, setzen sich zusammen aus der Suchanfrage, dem Zeitpunkt der Suche, den Userdaten, der angeklickten Position im Suchresultat und den Dokumentinformationen zum angeklickten Dokument.

### Wie können wir Click-Trough-Daten aus Webtrekk lesen?

Webtrekk ist ein Analysetool. Das heißt für uns, wir können nicht direkt auf die Datenbank mit den Tracking-Daten zugreifen. Um die Tracking-Daten lesen zu können, müssen wir eine Analyse auf Webtrekk ausführen. Mithilfe dieser Analyse können wir uns die Click-Trough-Daten so zusammenstellen lassen, wie wir sie für die Berechnung der Click-Trough-Rate benötigen.

Die Click-Trough-Daten bestehen aus einzelnen *Klick-Werten*. Ein Klick-Wert beschreibt die Anzahl der Klicks, die zu einer bestimmten Suchanfrage auf ein bestimmtes Dokument gemacht wurden und auf welcher Position im Suchresultat sich dieses Dokument dabei befunden hat. Die

Webtrekk-Analysen geben uns eine Sammlung von Klick-Werten zurück. Wir können bei diesen Analysen die Klick-Werte nach Suchbegriffen oder auch Suchtermen filtern und den Zeitraum mitgeben, in welchen die Suchanfragen durchgeführt wurden. Des weiteren gibt es die Möglichkeit, weitere Filter wie die Anzahl zurückzugebender Klick-Werte oder auch den „Login-Status<sup>1</sup> des Users“ zu setzen.

### Wie sehen die Click-Trough-Daten aus?

Eine Beispiel für einen Klick-Wert wie er von einer Webtrekk-Analyse ausgespielt wird, sieht wie folgt aus:

Click-Trough-Daten	Klick-Wert
searchresult-1.Course.chronische Dyspnoe bei Erwachsenen.10621768.chronische Dyspnoe	5

Hier die Aufschlüsselung der Click-Trough-Daten:

Position	Dokumenttyp	Titel	ID	Suchterm
searchresult-1	Course	chronische Dyspnoe bei Erwachsenen	10621768	chronische Dyspnoe

Die Click-Trough-Daten lassen sich wie folgt lesen. In diesem Beispiel haben die User mit der Suchanfrage „chronische Dyspnoe<sup>2</sup>“ gesucht. Dabei haben sie das Dokument mit der ID 10621768 angeklickt. Dieses hat sich dabei auf der Position eins der Suchresultate befunden. Es wurde insgesamt fünfmal angeklickt in der gesuchten Periode.

### Das Userverhalten bestimmt über die Relevanz eines Klicks

Zur Definition der Click-Trough-Daten gibt es verschiedene Ansätze. Wie in [Joachims] beschrieben, können anhand des Verhaltens der User vor während und nach dem Klick Schlussfolgerungen für die Relevanz des Klicks gemacht werden.

Wir bauen dabei auf den Erkenntnissen aus [Joachims] aus. In diesen Analysen wurde untersucht, wie aussagekräftig ein Klick während eines Suchvorgangs ist.

## 2.1.3 Result-Reranking mittels PBM Algorithmus

## 2.2 Zusammenfassung

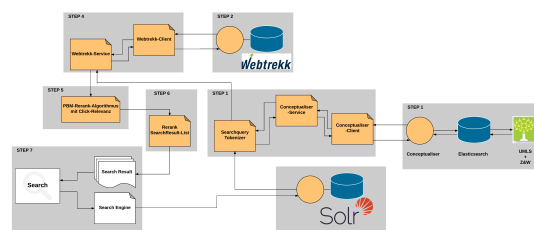
<sup>1</sup>Mit Login-Status wird zwischen einem zum Zeitpunkt der Suche auf der Springermedizin-Applikation angemeldeten und nicht angemeldeten User unterschieden

<sup>2</sup>Als Dyspnoe wird eine unangenehm erschwerte Atemtätigkeit bezeichnet

# Reranking mittels Click-Trough-Rate Ergebnis

## 3.1 Prozessaufbau des Lösungsansatzes

### 3.1.1 Prozessaufbau als Bild



**Abb. 4:** Prozessaufbau des Lösungsansatzes

### 3.1.2 Daraus entstehende Probleme

### 3.1.3 Suchterm Segmentierung

- Was wird gesucht und zu welchem Zweck Beispielsweise Zeitschrift => direkter Link auf Zeitschrift Oder Kurs (sehr tagesaktuell) => nicht Kern der Arbeit - Mehrdeutigkeiten - verwandte Begriffe und Synonyme - Segmentierung

- Contentprobleme: ————— - Mehrfachverwertung des Contents => Mehrfache Ausspielung der gleichen Artikel Grundsätzlich tritt dieses Problem bei allen Suchanfragen auf - z.B. auch bei den Anfragen COPD, Osteoporose, Sarkoidose und Schilddrüsenkarzinom. Für dieses Problem gibt es eine einfache Erklärung: News und kommentierte Studien laufen oft erst online und werden dann z.B. von der MMW, der PneumoNews und der InfoDiab zweitverwertet. Auf diese Weise taucht ein Text dann halt 2-4-6 mal auf ... Außerdem hätte ich gleich noch eine Frage. Bei den Suchanfragen wird mir zum Teil ein und derselbe Artikel gleich mehrfach präsentiert – was auch an der hier im Hause praktizierten Mehrfachverwertung liegt. Auf diese Weise sind dann vier von zehn Suchergebnissen identisch. Willst du das wirklich so haben?"

- Ausspielung von Teaser => Contentproblem Vielleicht noch eine Anmerkung zu „Teaser“ Ich denke, dass der Begriff jeden ernsthaften User dazu veranlasst, das nicht anzuklicken. Teaser sagt ja nichts über die Wertigkeit des Beitrages aus wie z.B. Original paper oder Review paper, Teaser sollte also wie Announcement gar nicht erscheinen. Aber darauf haben Sie wahrscheinlich keinen Einfluss.

- Habe alle Analysen nach der voreingestellten „Relevanz“ gemacht und nicht nach „Aktualität“. Dadurch wird das Ergebnis negativ beeinflusst, da eine Arbeit von 1998 heute keine Relevanz mehr hat, auch wenn der Begriff und die Art der Arbeit (Original paper) dafür sprechen. - Es wurden bei einer Reihe von Begriffen auch „Teaser“ gefunden. Da keine Autoren angegeben sind und man nicht weiß, ob entsprechend verknüpft ist, habe ich auch diese als nicht relevant eingestuft. - Manche Suchbegriffe sind aus meiner Sicht zu allgemein („operative Therapie“) oder veraltet („Prostata-Adenom“ statt benigne „Prostatahyperplasie“) - Eine weitere Unsicherheit entsteht durch den fehlerhaften Import der Daten und die falsche Rubrizierung im Rahmen der Überführung ins neue System – Vieles läuft unter Original Paper, das nicht korrekt ist. Auch der „article type“ sollte bei der Suche eine Rolle spielen. „original paper“ und „review paper“ und „continuing education“ vor „news“ und „report“ und: „announcement“ sollte gar nicht gefunden werden, oder?

### 3.1.4 Aufbereitung Click-Trough-Daten

#### Kein Einfluss auf Suchergebnisqualität während der Klicks

siehe Grundlagen

#### Userverhalten

Die Webtrekk-Analysen bieten uns 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.

### 3.1.5 Click-Trough-Rate in Suchprozess einbinden

- Solr gebundene Suchresultatmenge - Pagination (Einfluss von Reranking)

### 3.1.6 Result-Reranking mittels PBM Algorithmus

- Smoothing

- Mehrfachverwertung von Content => mehrfache Auflistung in Suchergebnissen (nicht Teil dieser Arbeit) - Algorithmus nicht festfahren (Overfitting)

- User nicht festfahren auf altem Wissen - aktuelle und neue Publikationen werden nicht berücksichtigt => keine Userrelevanz (nicht Teil dieser Arbeit)

## 3.2 Methodik

In Kapitel 2.1 haben wir gelernt wie Click-Trough-Daten entstehen und wie sie zu lesen sind. Nun können wir mit diesem Wissen die Click-Trough-Rate der Dokumente berechnen. Mithilfe der berechneten Click-Trough-Daten werden wir dann ein *Reranking* der Suchresultate durchführen. So wollen wir die Userrelevanz in die Suche einbinden. Die Vorgehensweise dazu sieht wie folgt aus.

## 3.2.1 Suchterm Segmentierung

### 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 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>1</sup> durchführen. Hier könnten wir uns überlegen, zusätzlich mit Stoppwörtern<sup>2</sup> nicht relevante Wörter aus dem Suchterm zu entfernen. Dieses Verfahren macht aber im Springermedizin-Kontext keinen Sinn. 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>3</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 [UMLS]) zugegriffen werden kann.

## 3.2.2 Aufbereitung Click-Trough-Daten

### Jeder Klick auf ein Dokument ist relevant

Wie in Kapitel 3.1.2 beschrieben, reichen Webtrekk-Analysen für komplexe Auswertungen der Click-Trough-Daten nicht aus. Wir können darum in dieser Arbeit *Feedback-Strategien* für die Click-Trough-Rate Auswertung, wie in [Joachims] beschrieben, nicht verwenden. Stattdessen müssen wir davon ausgehen, dass jeder Klick auf ein Dokument relevant ist.

### Gewichtung der Click-Trough-Daten

Durch die semantische Aufschlüsselung des Suchterms haben wir verschieden starke Relationen zwischen Click-Trough-Daten und dem Suchterm. Die Gewichtung der Stärke dieser Relation ist aber nicht Kern

---

<sup>1</sup>Bezeichnet die Aufteilung in Abschnitte, in diesem Fall in einzelne Worte

<sup>2</sup>Stoppwörter sind Wörter, die sehr häufig auftreten und für gewöhnlich keine Relevanz für den Dokumentinhalt besitzen

<sup>3</sup>Als Thesaurus wird ein strukturiertes Verzeichnis von Begriffen, welche allesamt in irgendeiner Beziehung stehen bezeichnet



dieser Arbeit. Wir gehen darum davon aus, dass unabhängig der Stärke der Relation zum Suchterm, alle Click-Trough-Daten eine gleiche Relevanz besitzen.

### Berechnung der Click-Trough-Rate

Die Click-Trough-Rate wird vor allem im Bereich des Internet-Marketing verwendet und stellt grundsätzlich die Anzahl der Klicks auf ein Dokument oder Link im Verhältnis zu den gesamten Impressionen dar. Gehen wir davon aus, könnten wir die Click-Trough-Daten eines Dokuments ins Verhältnis zu allen Click-Trough-Daten für einer Suchanfrage stellen und den Wert als Click-Trough-Rate verwenden. Wie wir aber bereits in Kapitel 2.1.1 gelernt haben, würden wir damit viele Problemstellungen der Interaktion der User mit der Suche ignorieren. Deswegen verwenden wir den Position-Based Modell basierten Algorithmus um die Click-Trough-Rate zu berechnen.

### 3.2.3 Click-Trough-Rate in Suchprozess einbinden

Es gibt drei mögliche Eingriffspunkte während des Suchprozesses, um die Click-Trough-Rate in der Springermedizin-Suche zu verwenden. Möglich wäre die Verwendung der Click-Trough-Rate in der Aufbereitung der Suchanfrage auf der Springermedizin-Applikation. Denkbar wäre auch, die Berechnung der Click-Trough-Rate in den Suchindex der Solr einzubauen. Die dritte Variante wäre die Verwendung der Click-Trough-Rate 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 Click-Trough-Rate direkt in die Solr einzubeziehen gibt es zwei Varianten. Wir können das Schema des Suchindexes über die Schema API (siehe [SchemaAPISolr]) erweitern und alle Einträge neu indexieren, oder wir ergänzen den Index um ein externes Feld (ExternalTextField) (siehe [ExtFieldSolr]).

Beide Lösungsansätze ergeben nur bei der Speicherung einer einfachen Click-Count Popularität<sup>4</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 [DisMax]). 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>5</sup>. Alle Boost-Faktoren hängen voneinander ab und müssten bei jeder Ergänzung um neue Faktoren normalisiert werden, um kein „über-Boosting“<sup>6</sup> einzelner Faktoren zu riskieren. Zudem besteht die Gefahr des „blinden

<sup>4</sup>Kennzahl für alle Klicks auf ein Dokument unabhängig des Suchterms

<sup>5</sup>Die Solr besitzt eine Boosting-Funktion, um bestimmte Wertübereinstimmungen in der Suche höher gewichtet zu können

<sup>6</sup>Bezeichnet die über-priorisierte Bewertung einzelner Faktoren

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 [**pbm**]) basierten Algorithmus zu verwenden. Dieser soll die Suchergebnisliste analysieren, die Click-Trough-Rate der Dokumente berechnen und die Liste neu sortieren.

Diese Lösung können wir relativ einfach in die Springermedizin-Applikation integrieren, ohne die restliche Suchlogik<sup>7</sup> zu beeinflussen. Hierbei müssen wir jedoch beachten, dass die Solr durch die Pagination-Funktion (siehe [**Pagination**]) 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.

### **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 Click-Trough-Rate-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.

## **3.2.4 Result-Reranking mittels PBM Algorithmus**

### **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 [**pbm**]). 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.

---

<sup>7</sup>Dazu gehört die Aufbereitung der Suchanfrage für die Solr und die Suche auf der Solr

## 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. Dies bestätigt die in Abbildung 3 dargestellte Analyse der Klicks auf die ersten 20 Positionen eines Suchergebnisses. Wir sollten darum darauf achten, dass je *schlechter* der Rang des angeklickten Dokumentes im Suchresultat der Solr ist, desto *höher* das Relevanzfeedback zu bewerten ist.

Das machen wir, indem wir das Verhältnis zwischen Klick-Wahrscheinlichkeit der Position und Klick-Wahrscheinlichkeit des Dokumentes abhängig der Position im Suchresultat definieren. In den ersten zehn Positionen des Suchresultat der Solr verwenden wir für die Berechnung der Click-Trough-Rate die beiden Klick-Wahrscheinlichkeiten mit gleicher Gewichtung. Für die Suchresultate zwischen 10 und eingeschlossen 20, setzen wir die Klick-Wahrscheinlichkeiten ins Verhältnis eins zu zwei, so dass die Klick-Wahrscheinlichkeit des Dokumentes eine stärkere Gewichtung erhält. Für die Suchresultate mit einer Position über 20 verstärken wir die Gewichtung der Klick-Wahrscheinlichkeit noch einmal und setzen die Wahrscheinlichkeiten ins Verhältnis eins zu drei. Das liegt daran, dass bei Klicks auf Dokumente mit einer solch hohen Position wir davon ausgehen können, dass die suchende Person die Suchresultate genau analysiert hat, bevor sie ein Dokument angeklickt hat. Haben wir die Verhältnisse definiert, müssen wir diese in den Algorithmus einbauen. Wie wir dies machen, werden wir im folgenden Abschnitt anschauen.

### 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* (siehe [weightedAVG]), dessen Gewichte sich zu Eins aufsummieren. Diese Gewichte sind die Smoothing-Faktoren, weshalb das Verfahren zählt zu den Smoothing-Algorithmen zählt.

### 3.2.5 Vergessen der alten Daten

Ein Algorithmus zur Berechnung von Wahrscheinlichkeiten muss sich ein gewisses Grundwissen aneignen. Dies geschieht üblicherweise durch Trainingsdaten. Genauso muss er alte Daten wieder vergessen können, um Overfitting<sup>8</sup> zu vermeiden.

---

<sup>8</sup>Überanpassung des Algorithmus durch zu viele (falsche oder veraltete) Daten

### Durch Webtrekk ist kein komplexer Lern-Algorithmus notwendig

Durch Webtrekk haben wir eine Wissensbasis, die sich stetig und zeitnah aktualisiert. So muss der Algorithmus nicht stetig neues Wissen lernen und altes vergessen, sondern er kann direkt diese Wissensbasis zugreifen. Dies geschieht, indem zur Laufzeit<sup>9</sup> Analysen gegen Webtrekk über eine frei definierbare Periode gemacht werden. Dadurch kann *Overfitting* vermieden werden. Deshalb verwenden wir keinen komplexen Lern-Algorithmen wie in [IWUSBI] vorgestellt.

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

Nun könnten wir die Klick-Wahrscheinlichkeit als absoluten Wert für die *Userrelevanz* betrachten. Dies wäre jedoch falsch, wie in Abbildung ?? dargestellt, müssen wir davon ausgehen, dass viele User der Qualität der Suchmaschine vertrauen. Diese betrachten die Top-Suchresultate als die relevanten Suchresultate. Denkbar wäre auch, dass User unabsichtlich das falsche Dokument anklicken und dadurch die Click-Trough-Rate eines Dokumentes verfälschen. Dadurch kann ein *Overfitting* des Algorithmus entstehen.

### Overfitting vermeiden

Um 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* arbeiten.

## 3.3 Der PBM-Algorithmus

## 3.4 Zusammenfassung

---

<sup>9</sup>Unter Laufzeit wird in diesem Fall der Zeitpunkt der direkte Abfrage während der Suchanfrage bezeichnet

# Implementierung

- 4.1 Technologie-Stack
- 4.2 Architektur der Implementierung
- 4.3 Highlight: Webtrekk-Analysen
- 4.4 Highlight: PBM Rerank-Algorithmus
- 4.5 Zusammenfassung

# Evaluation und Auswertung

## 5.1 Einführung

### 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.

## 5.2 Aufbau der Analyse

### 5.2.1 Datengrundlage

#### Filterung der nutzbaren Daten mittels Cohens Kappa

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 [Kappa]). 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.

### 5.2.2 Metrik

#### Evaluationsdaten mittels NDCG-Algorithmus auswerten

Um das Qualitätsmaß der beiden Suchen vergleichen zu können werden wir den Bewertungsalgorithmus *NDCG* (siehe [NDCG]) 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.

## **5.2.3 Vorgehen**

### **Evaluationssystem aufbauen**

Um eine Evaluation durchführen zu können, 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 zehn Suchergebnisse nach Relevanz zum Suchterm bewertet werden. Der Ergebnisse werden in einer Datenbank gespeichert, um sie später auszuwerten.

### **Evaluationssystem auswerten**

Nach Ablauf der Evaluationsphase werden wir die Evaluations-Daten auswerten. Die Auswertung der Daten findet direkt im Evaluationssystem statt.

Dazu werden die Daten aus der Datenbank gelesen und mit dem Cohens Kappa Koeffizienten die nutzbaren Daten gefiltert.

## **5.2.4 Durchführung**

### **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.

Hier folgend die Aufteilung der generierten Analysen:

- Springermedizin.de Suche
  - 50% aller Analysen
- Reranking Suche
  - CTR-Daten der angemeldeten User: 25% aller Analysen
    - \* Einfluss Zufallsranking 0.1: 50% dieser Analysen
    - \* Einfluss Zufallsranking 0.01: 50% dieser Analysen
  - CTR-Daten aller User: 25% aller Analysen
    - \* Einfluss Zufallsranking 0.1: 50% dieser Analysen
    - \* Einfluss Zufallsranking 0.01: 50% dieser Analysen

## 5.3 Auswertung der Suchergebnis-Qualität

### 5.3.1 Quantitative Auswertung

### 5.3.2 Diskussion

## 5.4 Zusammenfassung



# Zusammenfassung und Ausblick

## 6.1 Zusammenfassung

## 6.2 Ausblick



# Abbildungs-Verzeichnis

1	Aufbau der Suche bei Springer Nature . . . . .	2
2	Analyse der 20 am häufigsten angeklickten Dokumente der zehn meistgesuchten Suchphrasen. <i>Zeitraum der Analyse: 19.08.16 - 19.09.16</i> . . . . .	7
3	Analyse der Klicks auf die ersten 20 Positionen der Suchergebnisse aller Suchanfragen. <i>Zeitraum der Analyse: 19.08.16 - 19.09.16</i> . . . . .	8

## Tabellen-Verzeichnis