# Personalized Search

Building a prototype to infer the user's interest

### Tom Burgmans

Technology Product Owner Search Wolters Kluwer April 28<sup>th</sup> 2022







# Make search better through **Personalization**....a controversial topic





**Personalized search** is web based search results that are tailored specifically to an individual's interests by incorporating information about the individual beyond the specific query provided.



# What could make search "personal"





# What could make search "personal" (this PoC)





# Hypothesis

The basis for **Personalized Search** is a **Recommendation Engine** which is based on an index of **user activity** that allows to infer the **user's interest** via **collaborative filtering**.





# Collaborative filtering & anomaly detection

- Find users with **similar interests**
- Recommend items with an **unusually high presence** in this group compared to the rest





# Collaborative filtering & selecting the foreground

- Find only users with very similar interests
- Balance between similarity and volume





# What is the ideal time window to infer the user's interest?





# Example extremely temporary personalization





X 🤳 Q

🔍 All 🔚 Images 🖽 News 🖓 Maps 🗈 Videos 🗄 More Settings Tools

About 81.800.000 results (1,04 seconds)

en.wikipedia.org > wiki > John\_Lennon 👻

### John Lennon - Wikipedia

John Winston Ono Lennon MBE (born John Winston Lennon, 9 October 1940 – 8 December 1980) was an English singer, songwriter and peace activist who gained worldwide fame as the founder, co-lead vocalist, and rhythm guitarist of the Beates.

Other names: John Winston Ono Lennon Parent(s): Alfred Lennon, Julia Stanley

### Murder

lennon

On the evening of 8 December 1980, English musician John ... Mark David Chapman The Catcher in the Rye - Holden

Years active: 1956-1975, 1980

Died: 8 December 1980 (aged 40); New Yo ...

Caulfield - Chapter 27 - ...

solr × Q Q Q All 🕻 Images I Videos 🗐 News 🛇 Maps I More Settings Tools

About 19.000.000 results (0,62 seconds)

Ad · www.innoventsolutions.com/solr \*

### Solr Experts - InnoventSolutions.com

Solr Consulting, Training, Support. Implement, Configure, Optimize. Optimize Revenue. Request A Demo. Merchandise Search. Services: Solr installation, Configuration And Integration, Architecture For Disaster Recovery And Replication. Contact Us - Solr Advisory Services

Ad - www.opensolr.com/ \*

### Solr as a Service, Worldwide | Hosted Solr solutions for Free

Hosted Solr solutions in the cloud, Worldwide, High-Availability, Multi-Region Failover systems, all in one, fully self managed service. Solr Search Made Easy, Services: Shared Solr Cloud, Dedicated Solr Cloud. Enterrise Solr Clusters.

	<b>J</b>	
Q	luc	× 🏮 Q
HĮCY	Lucy 2014 film	Remove
0	lucene morelikethisquery	Remove
0	lucene boolean query isscoring	Remov
Ú.	Lucy In The Sky With Diamonds (Take 1 And Speech At The End) Song by The Beatles	
Q	lucy in the sky with diamonds lyrics	
n jezh	Lucardi	
į.	Lucifer Television series	
Q	lucifer season 5	
Q	lucid dreams lyrics	
Q	lucifer season 6	
		Report inappropriate prediction

Q	luc	× 🌷 🔍
HĮCY	Lucy 2014 film	Remove
0	lucene morelikethisquery	Remove
0	lucene boolean query isscoring	Remove
Q	lucene	
Q	lucene index	
Q	lucene search	
Q	Lucidworks Software company - Cambridge, UK	
Q	lucidchart	
100 Jacob	Lucardi	
ξt.,	Lucifer Television series	
		Report inappropriate prediction



# How to apply the inferred user's interest?

as a boost:



### as a filter:



### as interesting items:



### as a separate result set:





# Goals of this PoC

# How to build it?

What can we do on our existing technology stack?

# In what cases to apply it?

Where is the need for personalization? What form is needed?

# What data do we need?

Define 'additional' data sources; cleaning of noise

# How to measure success?

Can we test/tune it offline? How to A/B test online?



- Technology
- Need
- Data
- Verification







JLH score in Elastic = (fg\_percentage - bg\_percentage) \* (fg\_percentage / bg\_percentage)





Solr's SignificantTerms score = (log(freq\_foreground)+1) \* (log( (numDocs + 1)/(freq\_background + 1) ) + 1) (SignificantTermsStream.java)



# Finding a technology for anomaly detection





sqrt(fg size \* bg prob \* (1 - bg prob))



(RelatednessAgg.java)



# Selecting similar users

- MLT standard Solr query parser to find all similar documents
- MLTPLUS home made query parser to only keep documents that are *very* similar.

## Wrapping mlt into mltplus

```
(
({!wkmltplus minqt=3 mm=3 mmp=20 score=false}{!mlt v=$user maxqt=200 mintf=1 mindf=1 qf='suas'}) OR
({!wkmltplus minqt=10 mm=5 mmp=10 score=false}{!mlt v=$user maxqt=200 mintf=1 mindf=1 qf='docdocviews'}) OR
({!wkmltplus minqt=10 mm=10 mmp=70 score=false}{!mlt v=$user maxqt=200 mintf=1 mindf=1 qf='publicationdocviews'}) OR
({!wkmltplus minqt=10 mm=20 mmp=85 score=false}{!mlt v=$user maxqt=200 mintf=1 mindf=1 qf='palevel3'})
```

```
AND
```

```
({!filters param=$scope})
```





# Design recommendation engine

Source: Wolters Kluwer Navigator 2021 customer usage data.





DWH

# Personalization prototype









# What is the need for personalization based on a year history?

How consistent is the user's interest through the year?



# What is the need for personalization based on the last session ?

How much does the last session predict the interest of the next one?



For 42% of the users the previous session has a 0.9-1.0 consine overlap with the current one w.r.t. PA interest.

This rises to 66% when the time between sessions get shorter.

# How to measure success?

- Offline:
  - Manual tests by a domain expert
  - Build a model that predicts the likelihood that the previous session predicts the interest of the current one.
  - Replay user sessions and guestimate how much shorter they could have been in case personalization would have been applied.
- Online:
  - Measure clicks on explicit recommendations
  - A/B test (in case recommendations applied as boosts/filters)



# Lessons learned

### Yes we can

We can indeed infer the user's interest via collaborative filtering of user actions, using Solr and using usage data from a large Wolters Kluwer product

## What works best

Short-term personalization is applicable to **more users** than long-term personalization (may differ per business case)

## **Quality logs needed**

The accuracy of the recommendations depends on the **purity/completeness** of the **usage logs**. Preferred: front-end usage logging

### **Challenge: testing**

Offline testing of personalization is **hard**, but slightly easier for the shortterm variant. Online testing **takes time**.

## We're not done yet

Improvement opportunities:

- Use frond-end logs
- Better data cleaning
- Query normalization
- Add more signals like filters & autocomplete clicks
- Testing strategies

