

UNIVERSITY OF SOUTHERN DENMARK

SOFTWARE ENGINEERING 6. SEMESTER

---

## Datamining and its use

---

*Author:*

Lasse Bjørn HANSEN  
Simon FLENSTED

*Supervisor:*

Jan Corfixen Sørensen SØRENSEN



UNIVERSITY OF SOUTHERN DENMARK

*A report submitted in fulfillment of the  
requirements  
of Software Engineering 6. semester  
at*

University of Southern Denmark  
TEK

May 8, 2017

*“Some quote”*

- Gruppe 3

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.1.1	Data mining . . . . .	1
1.1.2	Product recommendation . . . . .	1
<b>2</b>	<b>Problem statement</b>	<b>2</b>
2.1	Problem description . . . . .	2
2.2	Problem statement . . . . .	2
<b>3</b>	<b>Implementation</b>	<b>3</b>
3.1	Solution overview . . . . .	3
3.2	Initial data dump . . . . .	3
3.3	Data Transformation . . . . .	4
3.4	The product recommendation algorithm . . . . .	6
3.5	Hosting the API . . . . .	8
<b>A</b>	<b>API commands</b>	<b>10</b>
	<b>Bibliography</b>	<b>12</b>



# Chapter 1

## Introduction

### 1.1 Motivation

The amount of data being processed around the Internet and within big systems is continuously increasing. This data should be structured and modelled in a way that makes it easily accessible and easy to work with. Handling large amounts of data the right way can prove to be very useful, not only to the company who possess the data, but also to the end users of a product. To achieve this, the art of data mining is very useful. The company Struct A/S have provided a typical software engineering task where data mining will create the foundation. This report will address theoretical aspects about data mining, how it is done in practice and how the final results of the processed data can be put to use. [1]

#### 1.1.1 Data mining

Data mining has become a big part of modern software engineering. Lots of companies tends to store large amount of data without structure and order within the data. This results in a lot of useless data which is both ineffective and a waste of resources. With prober data mining, it is possible to make this useless data useful to the company and its end users. In this case the data contain valuable information about users visiting the websites created and hosted by Struct A/S. By processing the data properly, it can be used for product recommendation, among other things.(Find kilde på dette)

#### 1.1.2 Product recommendation

If an e-commerce company wants to increase its profit, there is no doubt that product recommendation is one of the better ways of increasing your profit. This was first made popular by the retail giant Amazon. If you can predict what sorts of products your costumer may find useful, additional sales becomes more frequent. Big data sets, like the one provided by Struct A/S, can make it possible to predict customer needs, if the data is processed properly.(Find kilde på dette)

## Chapter 2

# Problem statement

### 2.1 Problem description

The initial problem/challenge is given to us by the company Struct A/S and is described as follows:

When launching sites, whether it being regular websites or web shops, a lot of user activity is logged. We therefore have a large amount of data associated with each of our sites but do not currently use it.

In the future we would like to be able to use logged data to generate an insight into the user activity on our site and actively use this data to create a personalized experience for the users.

This project will handle the initial normalization of the data, storing it in a scalable way and utilizing the data to create features which add value for the company and the users. In order to achieve this, theory has to become implementation. Research is required in terms of data storing, data mining and recommendation algorithms. This research is implemented in the end system creating an API allowing Struct A/S to get useful information from the data such as the recommended products for a certain user. This API will be the final product and will utilize different technologies and algorithms.

### 2.2 Problem statement

The data we have been given is in a de-normalized format and the problem therefore comes with two challenges - normalizing the data and utilizing the data to create a personalized experience for the users.

This leads to the following research questions:

- How can you effectively normalize large amounts of data?
- How can you optimally store and access data in a scalable way?
- How can you utilize the data to generate useful features for the company and the end user?

## Chapter 3

# Implementation

### 3.1 Solution overview

The final solution consists of a *RESTful* API build with ASP.NET Core and a MongoDB *No-SQL* database. Both are hosted through Amazon Webservices in an EC2 container using Docker. The API consists of the commands seen in appendix A.

This section covers the process from the initial data delivery to the final solution in a chronological order.

### 3.2 Initial data dump

In the beginning of the project we were supplied with data from one of *Struct A/S* customers. This data was in the form of many SQL tables and most of the data was not relevant for creating product recommendations. The main tables used were the following:

- Visitor: A collection of every unique visitor who visited the customer's website, each visitor gets a unique identifier called UID. Contains 3,073,665 visitors.
- Profile: A collection of every users signed up at the website. Contains 3037 profiles.
- BehaviorData: A collection of every unique action performed by visitors on the website, for example when a visitor views a product a new row is made with the visitor's UID, the product UID and the timestamp. Contains 3,326,736 visitor actions.
- Order: Contains each profile's orders. Contains 5520 orders.
- Product: A collection of all products on the website with their unique IDs. Contains 22,445 products.
- ProductGroup: A collection of all product groups. Contains 262 product groups.
- AttributeValueRendered: A collection of product and product group descriptions in different languages. Contains 409,259 descriptions:

A snapshot of the visitor and behaviorData table can be seen in figure 3.1 and 3.2.

	Uid	UserAgent	BrowserName	BrowserVersion	IPaddress
1	6820EDD0-E6A6-4105-A078-0000127E7AE1	AdsBot-Google (+http://www.google.com/adsbot.html)	Unknown	0.0	66.249.134.101
2	4D9FC023-9034-4B93-96D8-000013AB1940	Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.3...	Chrome	36.0	212.71.232.131
3	CC43266B-5E13-473E-A4F2-000019D2B065	Mozilla/5.0 (compatible; Googlebot/2.1; +http://ww...	Mozilla	0.0	66.249.94.101

FIGURE 3.1: Visitor table from the original data

	Type	Id	UserId	Timestamp	ActivityData
1	ProductView	31070	794ebe68-49c5-451c-ba80-226cdc0508f4	2016-09-30 20:04:01.5376346 +02:00	NULL
2	ProductView	31071	055bf3eb-a556-4672-9f65-5f64932e3973	2017-01-05 12:18:05.5065939 +01:00	NULL
3	ProductView	31071	06d3f1b4-8347-499d-9edd-71c058024cc1	2017-01-27 10:42:00.8031017 +01:00	NULL

FIGURE 3.2: Behavior data table from the original data

These tables contain all the pertinent information for creating product recommendations and can be utilized after a cleaning and structuring process. This process is described in the following sections.

### 3.3 Data Transformation

To begin the initial data transformation a data storing technology has to be selected. The technology chosen was *No-SQL*, specifically *MongoDB* the most popular No-SQL framework [4].

*No-SQL* is chosen because of the good fit for this project. The data demands are not clearly specified in the beginning and with No-SQL it is easy to add or remove data or even change the data types on the fly. No-SQL's denormalized format also allows for faster retrieval of a single item without having to do joins or complex SQL queries. Finally No-SQL is easier to scale across multiple servers and many engines have built in scaling functionalities [3] which can come in handy when multiple clients begin using the service.

A brief overview of the different terminology for SQL and No-SQL is given in table 3.1.

TABLE 3.1: SQL vs No-SQL terminology

SQL	No-SQL	Comment
Table	Collection	
Row	Document	A No-SQL document can contain more complex datatype compared to a row in SQL e.g arrays or other documents

Python was used to accomplish the early migration from SQL tables to MongoDB. Several scripts were created to retrieve the data from the SQL server and transfer it to the MongoDB database in the wanted format. Pseudo code of one of these scripts can be seen in algorithm 1.

After all the scripts are finished the No-SQL database has the collections Visitor and Product. A breakdown of documents in the two collections can be seen in table 3.2



**Algorithm 1** Product Script

---

```

1: SQLQuery = SELECT * FROM struct.Product
2: db = MongoDB
3: for all Rows r in SQLQuery do
4:     Product = {Id: r.id
5:         Description: ""
6:         Created: r.created
7:         visitorID: []
8:         ProductGroupId: 0 }
9:     db.insert(Product)
10: end for
11:
12: Description = ""
13: firstId = true
14: previousId = 0
15: previousGroupId = 0
16: SQLQuery = SELECT * FROM struct.product JOIN struct.attributeValueRendered ON id ORDER
    BY productId desc
17: for all Rows r in SQLQuery do
18:     currentId = r.ProductId
19:     currentGroupId = r.GroupId
20:     if firstId then
21:         previousId = currentId
22:         firstId = false
23:     end if
24:     if currentId != previousId then
25:         db.update({id: previousId},
26:             $set: db.description: description
27:             db.update({id: previousId},
28:                 $set: db.ProductGroupId: previousGroupId
29:                 description = ""
30:             end if
31:             description += r.description
32:             previousId = currentId
33:             previousgroupId = currentGroupId
34:         end for
35:         VisitorIds = []
36:         firstId = true
37:         previousId = 0
38:         SQLQuery = SELECT * FROM struct.BehaviorData
39:         for all Rows r in SQLQuery do
40:             currentId = r.Id
41:             if firstId then
42:                 previousId = currentId
43:                 firstId = false
44:             end if
45:             if currentId != previousId then
46:                 db.update({id: previousId},
47:                     $set: db.VisitorId: visitorIds
48:                     db.update({id: previousId},
49:                         visitorIds= []
50:                     end if
51:                     visitorIds.append(r.UserId)
52:                     previousId = currentId
53:                 end for

```

---

TABLE 3.2: An overview of the fields in each document in the collections

Document	Fields	Comment
Visitor	Id: string Behaviors: array ProfileUID: string CustomerUID: string	The behavior array is an array of documents with the fields Type, Id and Timestamp. This contains all the behaviors of the specific visitor
Product	Id: int ProductGroupId: int VisitorId: stringArray Description: string Created: DateTime	The visitorId array contains Ids of all visitors who have looked at this product

After the data has been cleaned and structured in No-SQL the algorithm for determining product recommendations can be made. The algorithm is described in the following section.

### 3.4 The product recommendation algorithm

There are multiple ways to implement a product recommendation algorithm all with their advantages and disadvantages. The method chosen for this project is called Item-to-Item collaborative filtering. Other methods and the reasoning why these weren't chosen is described in further detail in Chapter 6 discussion.

Item-to-Item collaborative filtering is a datamining tool to link items (products) with other items in terms of their similarity. This method is also the way *Amazon* handles their product recommendations [5].

The specifics of the algorithm differs from implementation to implementation. In this version each product is compared to other products based on how much they have been viewed together by customers, the likeness of their description and their product group.

The first run of the algorithm requires going through each product and the visitors of each product to see what else they have looked at. This needs a lot of resources, but once run only new behavior has to be re-calculated. A run down of the algorithm can be seen in algorithm 2.

**Algorithm 2** Item-to-Item collaborative filtering algorithm

---

```

for all Products p do
  productScores = Dictionary<int, double>
  for all Visitors in p do
    for all Products visitorProduct in Visitor behaviors do
      if productScores contains visitorProduct then
        productScores[visitorProduct]++
      else
        productScores.Add(visitorProduct, 1)
      end if
    end for
  end for
  Sort productScores after highest value
  for all Products similarProduct in productScores do
    productScores[similarProduct] = calculateSimilarityScore(p, similarProduct, productScores[similarProduct]) (see algorithm 3)
  end for
  Sort productScores after highest value
  Store top 10 productScores in database under p
end for

```

---

**Algorithm 3** Similarity calculations for two products

---

```

calculateSimilarityScore(mainProduct p1, compareProduct p2, currentScore)
  similarAttributeFactor = 0.02
  productGroupFactor = 0
  numOfSimAttributes = 0
  if p1.productGroup equals p2.productGroup then
    productGroupFactor = 0.02
  end if
  for all words w in p1.description do
    if w is in p2.description then
      numOfSimAttributes++
    end if
  end for
  if numOfSimAttributes equals 0 then
    similarAttributeFactor = 0
  end if
  return  $currentScore * (1 + productGroupFactor) * (1 + similarAttributeFactor^{numOfSimAttributes})$ 

```

---

After algorithms 2 and 3 each product in the database now has an array with the top 10 similar products based on amount of views, description and product group.

Next up is calculating the top products for each visitor, these are the products the specific visitor has viewed the most. This is accomplished by iterating through each visitor, checking their behavior and storing their top products as a field in the database. This calculation can be seen in algorithm 4. This calculation also requires a large amount of resources the first time, but very little to maintain.

---

**Algorithm 4** Calculations of each visitors top products

---

```

for all Visitors v do
  visitorProducts = Dictionary<string, int>
  for all Behaviors b in v do
    topVistorProducts[b.Id]++
  end for
  Sort visitorProducts after highest value
  Store top 5 visitorProducts in database under v
end for

```

---

Since all these calculations are made before the actual product recommendations are requested, the process of recommending products is quite fast. The recommendation process starts with retrieving the requested visitor's top products from the database, retrieving these products top similar products, sorting them based on their score and finally returning the amount asked for. The code for the recommendation part can be seen in algorithm 5

---

**Algorithm 5** Get product recommendations

---

```

visitorTopProducts = db.GetTopProducts(visitorUID)
productRecommendations = Dictionary<int, double>
for all Products p in visitorTopProducts do
  SimilarProducts = db.GetTopProductRecommendation(product)
  for all products simProduct in similarProducts do
    if productRecommendations contains simProduct then
      productRecommendations[simProduct] += similarProducts[simProduct]
    else
      productRecommendations.add(simProduct, similarProducts[simProduct])
    end if
  end for
end for
Sort productRecommendations after highest value
return amount of productRecommendations requested

```

---

The entire process of requesting product recommendation, running algorithm 5 and returning them takes less than 40ms which is one of the non-functional requirements.

Some other paths are required in certain situations such as when the visitor does not have any behavior or not enough behavior to satisfy the amount of recommendations requested. In these cases the remaining recommendations are filled from the top 20 most popular products in the last 30 days. The top 20 products are calculated by checking the timestamp and finding those in the last 30 days and then counting how many times each product was viewed. The top 20 products are stored in the database and can be calculated through an API call.

### 3.5 Hosting the API

An API can be hosted in several different ways, through many providers. This product recommendation API is hosted through Amazon WebServices in an *EC2 Instance* [6]. To accomplish this the ASP.NET core project is built in a Docker container, the container is pushed to the Docker Hub and then pulled and run in the *EC2 Instance*. The database is similarly packed in a docker container and

---

run in the *EC2 Instance*. The Docker containers have exposed ports to the rest of the internet and can be accessed via *EC2 Instance* public DNS or IP.



# Appendix A

## API commands

Function	URI	Example	Description
GET	recommendation/visitorUID/ numberOfRecommendations/ database	recommendation/AAF995AE-1DD0-41C6-898B-9CBEE884E553/5/Pandashop	Returns a JSON array of size numberOfRecommendations containing productUIDs which are the product recommendations for the specific visitor
PUT	visitor/visitorUID/database	visitor/AAF995AE-1DD0-41C6-898B-9CBEE884E553/Pandashop	Registers a new visitor with the database
PUT	product/productUID/description/productGroup/database	product/5352/Agreatproduct/5/Pandashop	Registers a new product along with its description and product group with the database
PUT	behavior/visitorUID/behaviorType/ItemID/database	behavior/AAF995AE-1DD0-41C6-898B-9CBEE884E553/ProductView/5352/Pandashop	Registers a new behavior for the specific visitor with the database
GET	Update/database/password	Update/pandashop/supersecretpassword	Builds the collaborative filter for the database
GET	Updatevisitorproducts/database/password	Updatevisitorproducts/pandashop/supersecretpassword	Updates the top products for all visitors
GET	calculateTop20/database/password	calculateTop20/pandashop/supersecretpassword	calculates the top 20 products in the last 30 days

# Bibliography

- [1] A. Trivedi. (Feb. 2014). Mapping relational databases and sql to mongodb, [Online]. Available: <https://code.tutsplus.com/articles/mapping-relational-databases-and-sql-to-mongodb--net-35650> (visited on 02/23/2017).
- [2] M. D. C. Bowen. (Jan. 2017). What is normalized vs. denormalized data?, [Online]. Available: <https://www.quora.com/What-is-normalized-vs-denormalized-data> (visited on 02/23/2017).
- [3] C. Buckler. (Sep. 2015). Sql vs nosql: The differences, [Online]. Available: <https://www.sitepoint.com/sql-vs-nosql-differences/> (visited on 02/23/2017).
- [4] S. IT. (May 2017). Db-engines ranking, [Online]. Available: <https://db-engines.com/en/ranking> (visited on 05/07/2017).
- [5] G. L. B. Smith and J. York. (Feb. 2003). Amazon.com recommendations: Item-to-item collaborative filtering, [Online]. Available: <https://www.cs.umd.edu/~samir/498/Amazon-Recommendations.pdf> (visited on 05/08/2017).
- [6] Amazon. (May 2017). Getting started with amazon ec2 linux instances, [Online]. Available: [http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/EC2\\_GetStarted.html](http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/EC2_GetStarted.html) (visited on 05/08/2017).