

Presentation of paper "Learning Multi-dimensional Indexes"

Christina Borovilou

ds1200008@di.uoa.gr

University of Athens, Department of Informatics and Telecommunications

Athens, Greece



Figure 1: Indexing before computers' era

ABSTRACT

As computer science becomes more sophisticated, hardware more powerful, data more massive and users more demanding, the need for strong and fast search is loud. We are in an era that we should exceed the capabilities of our tools, exploit logical relations among data in order to save unnecessary computational effort and to build strong systems to serve people's needs in every aspect. Searching for data is a fundamental practice that needs improvement in parallel with the technological evolution.

CCS CONCEPTS

• **Theory of computation** → Data structures and algorithms for data management; • **Information systems** → Database query processing; Data analytics.

KEYWORDS

Datasets, learning index, Multi-dimensional Index, MIT, Machine Learning

ACM Reference Format:

Christina Borovilou. 2021. Presentation of paper "Learning Multi-dimensional Indexes". In *Proceedings of Data Management and Applications Workshop (DMAW 2021)*. ACM, New York, NY, USA, 3 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

This work is a part of an educational project under the course "Database Systems" of Master Program DSIT of NKUA. The information discussed below is triggered by a recent paper of Vikram Nathan, Jialin Ding, Mohammad Alizadeh and Tim Kraska concerning Learning Multi-dimensional Indexes [?]. In this presentation, we will give an overview of the things the paper is proposing, procedures that reevaluates and ideas that are generated. We will demonstrate some ideas discussed at the paper that we implemented in the scope of the course project. We will showcase the ideas behind, how they worked and we will point out observations as well as pain points of the attempt.

2 PAPER OVERVIEW

The team in this paper after presenting the current "tools in Indexing industry", introduces a new index, named "Flood". Starting from the description of the new index, "Learning, Multi-dimensional Index" we distinguish two equal complicated characteristics of the index: it is multi-dimensional and also seems to commit training in order to learn and adapt. Thus, we will examine indexing from two different perspectives:

Unpublished working draft. Not for distribution.
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted by ACM, provided that the copies are not made for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
DMAW 2021, September 24, 2021, Athens, GR
© 2021 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

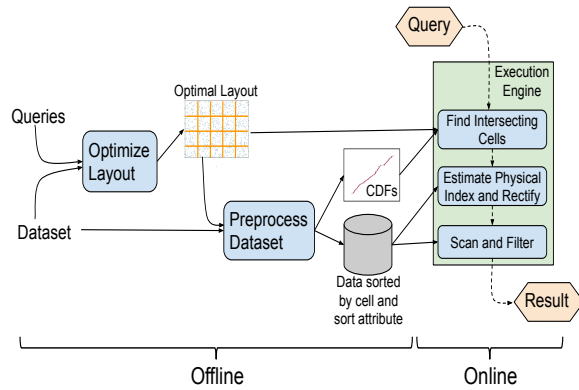


Figure 2: Flood's Architecture

Learning Index: They introduce the learning index to be the result of 2 important parameters: Learning from user's queries (frequency, ranges and correlations) as well as learning from data itself, in the sense that instead of adapting our Queries and optimizers to computer architecture, we can adapt our data architecture to the characteristics of our Database. As they notice though experiments to databases with different characteristics (Skewness, Size, Dimensionality, correlated attributes) query time parameters vary lot and do not follow a pattern that can be modeled analytically. And when math do not serve the problem, Machine Learning is called to find the parameters that make the system faster. To do so, Flood's team supports that query time is mainly the combination of 3 parameters according to its architecture:

- Projection: Identify the cells in the grid layout (we will explain later on) that mach the filter of SELECT statement
- Refinement: responsible to decide which attributes should be ordered and in which sequence
- Scan: time depending on users's number of filters, data distribution

These parameters do not contribute equally to the total time. Floods team uses ML algorithms to calculate these hidden weights. Flood trains a random forest regression model to predict them based on the output statistics. At this point they reach an interesting finding: The weight models are accurate across different datasets and query workloads! Independently of new data insertion or changes over query distribution, Flood will only reevaluate the existing models, will not restart the leaning process from the start!

Multi-Dimensional Index: For sure we cannot order data in more than one dimensions; but we can find heuristic ways to simulate it. We all know the value of binary search in query response. The core of that heuristic technique is that it imposes relational logic among data and exploits it when it is time. So if the one who "serves" data is the one who is storing data (meaning Database System) many techniques can be developed to the favor of easier access. That is what the Flood team introduces: a "geometrical" storage trick that help queries to feel more lucky even when index is on the filter.

```
SELECT SUM(R.X)
```

```
FROM MyTable
```

```
WHERE (a ≤ R.Y ≤ b) AND (c ≤ R.Z ≤ d)
```

Flood receives as input a filter predicate consisting of ranges over one or more attributes, joined by ANDs. The intersection of these ranges defines a hyper-rectangle, and Flood's goal is to find and process exactly the points within this hyperrectangle (e.g., by aggregating them). This feature we will examine in the Implementation Section.

Note that equality predicates of the form $R.Z = f$ can be rewritten as $f \leq R.Z \leq f$. We should mention that Floods work is focused on "ANDs" operations on WHERE clause. Equality predicates of the form $R.Z = f$ can be rewritten as $f \leq R.Z \leq f$. Disjunctions (i.e. OR clauses) can be decomposed into multiple queries over disjoint attribute ranges.

3 IMPLEMENTATION - OPTIMIZING THE GRID

Since the paper did not provide in public the code, my implementation was not to expand the work that was done but try to reproduce some part of it. In the following steps i will demonstrate the creation a new layout, adapted to specific attributes (Sequence of the attributes DOES matter). Sequence of attributes for the creation of layout in the scope of paper is given by ML procedures in the first part. In the code of the implementation is arbitrary.

3.1 Grid Layout Density

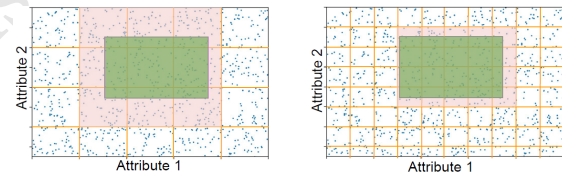


Figure 3: Scan Overhead

Layout is visualized above and is the grid that 2 (or more in the N-dimensional space) attributes scatter their data. The 2 representations above show off the query optimization dilemma: Split the space with dense grid or to a more sparse one? The first will achieve low scan overhead (it is the red area of the grid and corresponds to scanned points that didn't meet the selection criteria). but that would result to many visits (time rises up) in each cell and will counterbalance the profit of low scan overhead. The minimal loss will be given if we ask data. That process is in the ML process we mentioned in Section 2.

3.2 Data Distribution Skewness

Flood is also concerned about equally distributing data points along the layout. Using pandas in python we process our data creating equally distributed data in the grid and not a homogeneously in space (that will not offer any added value)

That feature will be used in the next section.

