# Report

## Comparison of BUC and AOI

### a. The primary purposes and use cases of each technique

Bottom-Up Cube (BUC) Algorithm

Purpose: Efficiently computes data cubes for OLAP by breaking down aggregate data into smaller sub-problems.

Use Case: Ideal for multidimensional analysis of large datasets, especially when some dimensions don't need full computation.

Attribute-Oriented Induction (AOI)

Purpose: A data mining technique for concept generalization, simplifying data by moving from lower to higher-level abstractions.

Use Case: Used in data summarization, reducing raw data to higher-level concepts.

### b. The types of insights or patterns each method is best suited to discover

Bottom-Up Cube (BUC) Algorithm:

**Insights**: Multidimensional aggregation, outliers in aggregates, and comparative patterns.

**Best Suited For**: Summarizing large datasets, business intelligence, and exploratory data analysis.

Attribute-Oriented Induction (AOI):

**Insights**: Concept generalization, hierarchical relationships, and data trend summarization.

**Best Suited For**: Knowledge discovery, data generalization, and summarization of broad patterns.

### c. The computational efficiency and scalability of each approach.

Bottom-Up Cube (BUC) Algorithm:

**Efficiency**: Optimizes cube computation by pruning irrelevant sub-cubes. Handles sparse data well but struggles with high-dimensional datasets.

**Scalability**: Performs better on sparse, low-dimensional data; suffers from the curse of dimensionality.

Attribute-Oriented Induction (AOI):

**Efficiency**: Reduces granularity, simplifying data through generalization. Efficient for structured data with defined hierarchies.

**Scalability**: Scales well for structured data but may struggle with complex hierarchies or simultaneous generalization of many attributes.

## d. The interpretability of the results produced by each method.

Bottom-Up Cube (BUC) Algorithm:

**Interpretability**: Highly interpretable, producing clear aggregated summaries, but high-dimensional results can be overwhelming.

**Ease of Interpretation**: Works well with business intelligence tools for interactive, multidimensional exploration.

Attribute-Oriented Induction (AOI):

**Interpretability**: Produces easily understood generalized concepts, but risks losing important details.

**Ease of Interpretation**: Simplifies complex datasets into digestible categories, ideal for high-level overviews.

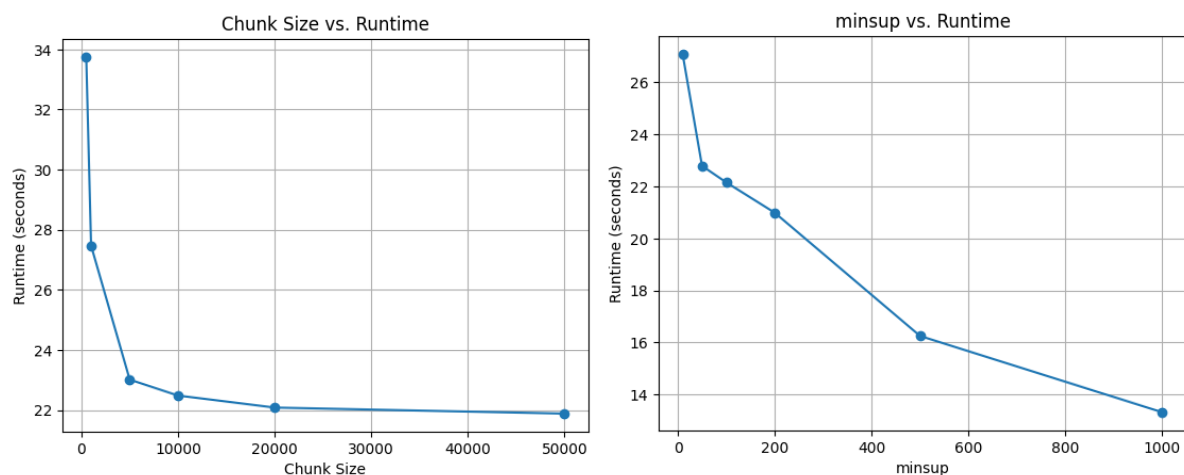## e. Scenarios where one method might be preferable over the other

Bottom-Up Cube (BUC) Algorithm:

1. **Multidimensional Analysis**: Ideal for analysing sales or transaction data across multiple dimensions.

2. **OLAP Decision Support**: Useful for institutions needing complex multidimensional queries.

3. **Sparse Data**: Optimized for scenarios with sparse datasets.

Attribute-Oriented Induction (AOI):

1. **Data Summarization**: Great for generalizing purchase behaviours or segmenting customers.

2. **Knowledge Discovery**: Suitable for identifying trends in large databases, such as healthcare or education.

3. **Concept Hierarchies**: Works well for summarizing performance or behaviour into higher-level categories.

From our implementation we got that AOI can self-decide and remove attributes, if necessary, but BUC cannot do self-decision. We need domain knowledge in case of AOI to create hierarchy table, and BUC uses domain knowledge to decide necessary and unnecessary attributes. AOI identifies generalised patterns in the data whereas BUC identifies most significant associations/correlations. In general BUC cannot handle large number of attributes as it must do many recursive calls which is not happening in AOI. So, AOI is much faster and more scalable than BUC. AOI results can be much more understandable to people doesn't know the data set also as the results are much summarized. If we want general customers behaviour AOI is sufficient but if we want any correlations between factors of that effect customer behaviour BUC is useful.



If the chunk size is too small, the system frequently loads data from disk. Disk access (I/O operations) is much slower than memory access, leading to significant delays in runtime. When you increase the chunk size, more data can be processed in memory at once, reducing the number of times the system needs to load new data from the disk. This reduces the overall runtime, as the system spends less time on I/O operations and more on actual computation.

Larger chunks allow more efficient use of memory and reduce the time spent switching between disk and memory. However, there is a limit—if the chunk size becomes too large and exceeds available memory, it can lead to memory swapping, which will negatively impact performance.

As you increase the minsup threshold, the number of patterns that meet the criteria decreases. This means that many combinations of dimensions or data points that don't meet the higher minsup threshold are "pruned" (ignored) by the algorithm. Fewer combinations mean fewer calculations and faster runtime, as the algorithm can skip large portions of the dataset.

When minsup is set too low, the algorithm spends time computing many patterns, including many that may be of little significance. By raising minsup, BUC can focus only on the most frequent and relevant patterns, significantly improving efficiency.

## Optimisation:

We applied counting sort for the BUC algorithm and ordered dimensions based on cardinality, skewness, and information gain, following our intuition that the most discriminating dimensions should come first. By default, we used counting sort in the algorithm with minimum support. Additionally, we implemented partitioning with external sorting and paging for cases where data doesn't fit in memory. As expected, out-of-memory scenarios like paging and partitioning with external sorting significantly increased processing time due to the overhead of disk access.

In terms of performance, we observed a runtime of 0.8s with no ordering, 0.72s when ordering by cardinality, and 0.66s for both skewness and information gain.