### **International Journal of Computer Engineering and Technology (IJCET)**

Volume 15, Issue 4, July-Aug 2024, pp. 516-526, Article ID: IJCET\_15\_04\_045 Available online at https://iaeme.com/Home/issue/IJCET?Volume=15&Issue=3

ISSN Print: 0976-6367 and ISSN Online: 0976-6375

Impact Factor (2024): 18.59 (Based on Google Scholar Citation)

DOI: https://doi.org/10.5281/zenodo.13310800





© IAEME Publication

# OPTIMIZING FUNNEL ANALYSIS IN MODERN DATA WAREHOUSES

## Satyam Shekhar

NetSpring Data Inc., USA



#### **ABSTRACT**

This article explores the implementation of funnel analysis in modern data warehouses, focusing on its importance for product managers in understanding and optimizing user journeys. It delves into the mechanics of funnel analysis, discussing two primary approaches: the Join Sequence and Stacked Window Functions methods. The article examines various query optimization techniques modern data warehouses employ, including common subexpression elimination, aggregate pushdown, and efficient handling of window functions. Additionally, it addresses performance considerations for both approaches, highlighting the benefits of pre-computed join indices and table clustering.

Throughout, the article emphasizes the critical role of funnel analysis in driving data-driven decision-making and product success in today's competitive business landscape.

**Keywords:** Funnel Analysis, Data Warehouses, Query Optimization, User Journey, Product Performance

**Cite this Article:** Satyam Shekhar, Optimizing Funnel Analysis in Modern Data Warehouses, *International Journal of Computer Engineering and Technology (IJCET)*, 15(4), 2024, pp. 516-526.

https://iaeme.com/MasterAdmin/Journal\_uploads/IJCET/VOLUME\_15\_ISSUE\_4/IJCET\_15\_04\_045.pdf

#### INTRODUCTION

Funnel analysis has become an indispensable tool for product managers seeking to understand and optimize user journeys through various stages of product interaction. This analytical approach allows managers to track and visualize the customer's progression from initial acquisition through activation, retention, referral, and ultimately revenue generation [1]. By mapping out these critical touchpoints, product managers can identify bottlenecks, optimize conversion rates, and make data-driven decisions to enhance overall user experience and product performance.

The power of funnel analysis lies in its ability to provide granular insights into user behavior at each stage of interaction. For instance, in an e-commerce platform, a typical funnel might include visiting the homepage, browsing products, adding items to the cart, initiating checkout, and completing a purchase. By examining the conversion rates and drop-offs between these stages, product managers can pinpoint areas for improvement and allocate resources effectively to maximize user engagement and conversions [2].

As the volume of user data grows exponentially, the challenge lies in collecting this information and efficiently processing and analyzing it to derive actionable insights. This is where modern data warehouses come into play, offering robust solutions for executing complex funnel queries at scale. These advanced systems are designed to handle massive datasets and perform sophisticated analyses with impressive speed and accuracy.

This article delves into the mechanics of funnel analysis, exploring its implementation within the context of modern data warehouses. We will examine different approaches to constructing funnel queries, discuss optimization techniques, and highlight how these powerful analytical tools can be leveraged to drive product success in today's data-driven business landscape.

## **Understanding Funnel Analysis:**

Funnel analysis is a powerful technique that tracks user progression through defined stages, providing product managers with a visual representation of conversion and drop-off rates at each step of the user journey. This method derives its name from the funnel-like shape that typically emerges when visualizing user flow, with the number of users decreasing as they move through successive stages [3].

At its core, funnel analysis helps product managers answer critical questions about user behavior and product performance. For instance, in an online learning platform, a typical funnel might include stages such as "Getting Started," "Analytics Course," and "Paid Activation." By examining this funnel, product managers can determine:

- 1. How many users complete the initial onboarding process?
- 2. What percentage of users progress from introductory content to more advanced courses?
- 3. What proportion of users ultimately convert to paid subscribers?

The insights gained from funnel analysis can be transformative for product strategy. For example, if there's a significant drop-off between "Getting Started" and "Analytics Course," it might indicate that users are struggling to find relevant courses or that the transition between introductory and advanced content needs improvement. Similarly, a low conversion rate to "Paid Activation" could suggest that the perceived value of premium content needs to be enhanced or that the pricing strategy should be reevaluated.

Moreover, funnel analysis enables product managers to:

- Identify bottlenecks in the user journey: By pinpointing stages where users are most likely to drop off, product teams can focus their efforts on improving these critical points.
- Prioritize feature development: Understanding where users struggle or lose interest helps in allocating resources to the most impactful improvements.
- Measure the impact of product changes: By comparing funnel metrics before and after implementing changes, teams can quantify the effectiveness of their interventions.
- Set realistic conversion goals: Historical funnel data provides a baseline for setting achievable targets and benchmarking performance.
- Segment users: Analyzing how different user cohorts move through the funnel can reveal valuable insights about user preferences and behaviors.

The power of funnel analysis lies in its ability to provide a clear, data-driven picture of the user journey. This visualization makes it easier for stakeholders across the organization to understand and act on user behavior patterns. For instance, marketing teams can use funnel insights to refine their messaging and target their efforts more effectively, while development teams can prioritize features that address the most significant drop-off points [4].

Furthermore, funnel analysis can be extended beyond the traditional product usage scenario. It can be applied to various business processes, such as sales pipelines, customer support workflows, or even employee onboarding procedures. This versatility makes funnel analysis a valuable tool across multiple departments within an organization.

By leveraging funnel analysis, product managers can make data-driven decisions to optimize each stage of the user journey, ultimately leading to improved user engagement, higher retention rates, and increased revenue. As businesses continue to prioritize user experience and customer satisfaction, the role of funnel analysis in shaping product strategy is likely to become even more crucial.

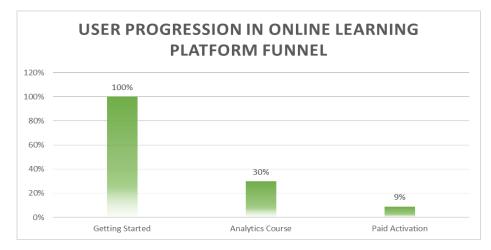


Fig. 1: Conversion Rates Across Key Stages of User Journey [3, 4]

#### **Data Model:**

To implement funnel analysis in modern data warehouses, we begin with a simplified schema centered around a single "Events" table. This table typically contains essential columns such as user\_id, event\_name, and event\_time. While this model makes some simplifying assumptions, it serves as a foundational structure for understanding funnel analysis implementation. It's important to note that modern data warehouses can handle far more complex scenarios efficiently, including semi-structured data and multiple table joins [5].

The Events table schema might look like this:

```
CREATE TABLE Events (
    event_id INT PRIMARY KEY,
    user_id INT,
    event_name VARCHAR(255),
    event_time TIMESTAMP,
    additional_attributes JSON
);
```

This schema allows for flexibility in capturing various event types while maintaining a structured core. The additional\_attributes JSON column can store event-specific data, accommodating the diverse nature of user interactions without requiring schema changes for each new event type.

When it comes to executing funnel queries, two primary approaches have emerged, each with its own set of advantages and challenges:

# 1. Join Sequence Approach:

This method relies on creating subqueries for each funnel stage and then joining the results to calculate user progression. The process involves:

- Generating a subquery for each stage of the funnel
- Joining these subqueries to track user movement through the stages

A simplified example of this approach might look like:

```
WITH stage1 AS (
  SELECT DISTINCT user_id
  FROM Events
  WHERE event_name = 'Getting Started'
),
stage2 AS (
  SELECT DISTINCT e.user_id
  FROM Events e
  JOIN stage1 s ON e.user id = s.user id
  WHERE e.event name = 'Analytics Course'
  AND e.event_time > (SELECT event_time FROM Events WHERE user_id = e.user_id AND
  event_name = 'Getting Started' LIMIT 1)
-- Additional stages...
SELECT
  COUNT(DISTINCT stage1.user_id) AS stage1_count,
  COUNT(DISTINCT stage2.user id) AS stage2 count,
  -- Additional stage counts...
FROM stage1
LEFT JOIN stage2 ON stage1.user_id = stage2.user_id
-- Additional joins...
```

## **Advantages:**

- Straightforward SQL syntax, making it easier for less experienced analysts to understand and modify
- Intuitive representation of the funnel stages in the query structure

## **Challenges:**

- Requires multiple table scans, which can be resource-intensive for large datasets
- Potentially expensive joins between subqueries, especially with high-cardinality user IDs

# 2. Stacked Window Functions Approach:

This more advanced technique leverages SQL window functions to create a stack of operations, one for each funnel stage. The process involves:

- Utilizing window functions to partition data by user and order by event time
- Creating a cumulative stack of these functions to represent each funnel stage

An example of this approach:

```
WITH funnel_stages AS (
 SELECT
   user_id,
    MAX(CASE WHEN event_name = 'Getting Started' THEN event_time END) OVER
 (PARTITION BY user_id) AS stage1_time,
    MAX(CASE WHEN event name = 'Analytics Course' THEN event time END) OVER
 (PARTITION BY user_id) AS stage2_time,
    -- Additional stages...
 FROM Events
SELECT
 COUNT(DISTINCT user id) AS total users,
 COUNT(DISTINCT CASE WHEN stage1 time IS NOT NULL THEN user id END) AS
 stage1 count,
 COUNT(DISTINCT CASE WHEN stage2_time IS NOT NULL THEN user_id END) AS
 stage2 count,
 -- Additional stage counts...
FROM funnel_stages;
```

#### **Advantages:**

- Performs a single table scan, significantly reducing I/O operations
- Results in a simpler relational plan, often leading to better query performance

#### **Potential optimization:**

• Clustering the Events table by event\_time can further enhance performance by reducing or eliminating the need for sorting operations

The choice between these approaches often depends on factors such as the analysis's specific requirements, the dataset's size, and the capabilities of the data warehouse in use. While the Join Sequence Approach might be more intuitive for simple funnels, the Stacked Window Functions Approach generally offers superior performance for complex, multi-stage funnels on large datasets [6].

It's worth noting that modern data warehouses continue to evolve, with many now offering advanced optimizations that can significantly improve the performance of both approaches. These optimizations may include intelligent query planning, automatic indexing, and adaptive execution strategies. As such, data analysts and engineers must stay informed about the latest features and best practices specific to their chosen data warehouse platform.

Approach	Table Scans	Query Complexity	Performance on Large Datasets	Ease of Understanding
Join Sequence	Multiple	High	Lower	Higher
Stacked Window Functions	Single	Medium	Higher	Lower

**Table 1:** Comparison of Funnel Analysis Approaches in Data Warehouses [5, 6]

# **Query Optimization Techniques:**

Modern data warehouses employ sophisticated optimization techniques to enhance the performance of complex queries, including those used in funnel analysis. These optimizations are crucial for handling large-scale data efficiently and providing timely insights. Let's explore some of the key optimization techniques:

# 1. Factoring out redundant computations:

Advanced query optimizers can identify and eliminate redundant calculations within a query. This process, known as common subexpression elimination, involves recognizing repeated subqueries or expressions and computing them only once. The results are then reused across the query, significantly reducing processing time and resource consumption [7].

For example, in a funnel analysis query, we might have:

```
SELECT

COUNT(DISTINCT CASE WHEN event_name = 'Sign Up' THEN user_id END) AS signups,

COUNT(DISTINCT CASE WHEN event_name = 'Purchase' THEN user_id END) AS purchases,

COUNT(DISTINCT CASE WHEN event_name = 'Sign Up' THEN user_id END) /

COUNT(DISTINCT user_id) AS signup_rate

FROM events

WHERE date >= '2023-01-01'
```

The optimizer would recognize that the COUNT(DISTINCT CASE WHEN event\_name = 'Sign Up' THEN user\_id END) expression is repeated and compute it only once.

## 2. Pushing aggregations below joins:

This technique, also known as aggregate pushdown, involves moving aggregation operations before join operations in the query execution plan. By performing aggregations earlier in the process, the amount of data that needs to be joined is reduced, leading to improved query performance. This is particularly beneficial in funnel analysis queries where aggregations on user actions are common before joining with other stages of the funnel.

## 3. Single sort across multiple window functions:

When dealing with multiple window functions that share the same partitioning and ordering criteria, modern optimizers can perform a single sort operation that serves all these functions. This is especially relevant in the Stacked Window Functions approach to funnel analysis, where multiple window functions are used to represent different stages of the funnel. By avoiding redundant sorting operations, query execution time can be significantly reduced.

For instance, in a query like:

SELECT user\_id,

MAX(CASE WHEN event\_name = 'Sign Up' THEN event\_time END) OVER (PARTITION BY user\_id) AS signup\_time,

MAX(CASE WHEN event\_name = 'Purchase' THEN event\_time END) OVER (PARTITION BY user\_id) AS purchase\_time

FROM events

The optimizer would perform a single sort operation on (user\_id, event\_time) to serve both window functions.

# 4. Local distinct aggregation within partitions:

For queries involving distinct aggregations, such as counting unique users in each funnel stage, optimizers can perform these aggregations locally within data partitions. This approach is particularly effective when data is already partitioned by the aggregation key (e.g., user\_id). By avoiding the need to shuffle data across nodes in a distributed system, this optimization can lead to substantial performance improvements [8].

These optimization techniques work in concert to dramatically improve query performance, especially for complex analytical queries like those used in funnel analysis. However, it's important to note that the effectiveness of these optimizations can vary depending on the specific data warehouse system, the nature of the data, and the complexity of the query.

Moreover, query optimization is an ongoing field of research and development. Data warehouse vendors continually refine their optimization strategies and introduce new techniques. For instance, some systems now employ machine learning algorithms to predict the most efficient query execution plans based on historical performance data.

To fully leverage these optimization techniques, it's crucial for data engineers and analysts to:

- 1. Understand the specific optimization capabilities of their chosen data warehouse system
- 2. Design schemas and write queries that can take advantage of these optimizations
- 3. Regularly review and update their approach as new optimization features become available

For example, when designing a schema for funnel analysis, consider:

- Partitioning strategies that align with common query patterns
- Appropriate indexing to support frequent join and filter operations
- Denormalization techniques to reduce the need for complex joins

By staying informed about these advanced optimization techniques and applying them judiciously, organizations can ensure that their funnel analysis queries - and indeed all their analytical workloads - are executed with maximum efficiency, enabling faster insights and better decision-making.

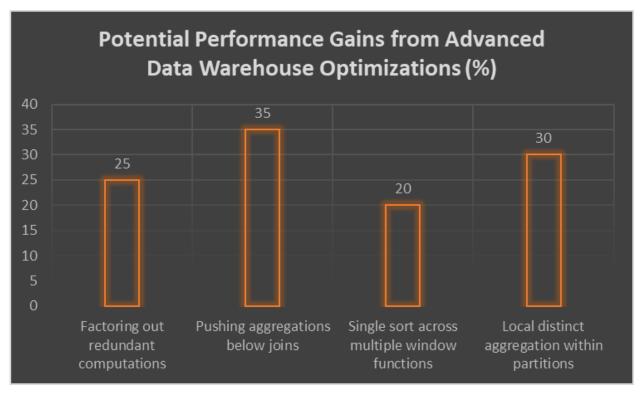


Fig. 2: Comparative Impact of Query Optimization Techniques in Funnel Analysis [7, 8]

## PERFORMANCE CONSIDERATIONS:

When implementing funnel analysis in modern data warehouses, it's crucial to consider the performance implications of different query approaches. Both the Join Sequence and Stacked Window Functions methods have unique performance characteristics that can be optimized further:

# 1. Join Sequence Approach:

The Join Sequence method, while intuitive, can face performance challenges, especially with large datasets. However, it may benefit significantly from pre-computed join indices. These indices are data structures that store the results of join operations, allowing for faster retrieval during query execution [9].

Pre-computed join indices can dramatically improve query performance by:

- Reducing the need for on-the-fly join computations
- Minimizing data movement across the cluster
- Enabling faster data access patterns

For example, in a funnel analysis scenario, we might create join indices that precompute the relationships between user actions across different funnel stages. This could involve creating an index that maps user IDs to their progression through the funnel stages, allowing for rapid retrieval during query execution.

However, it's important to note that while pre-computed join indices can offer substantial performance benefits, they also come with trade-offs:

- Increased storage requirements to maintain the indices
- Additional computational overhead to keep indices up-to-date as new data arrives
- Potential impact on data ingestion performance

Therefore, the decision to implement pre-computed join indices should be based on a careful analysis of query patterns, data volumes, and available resources.

# 2. Stacked Window Functions Approach:

The Stacked Window Functions method can leverage table clustering for improved performance. Table clustering involves physically organizing table data based on specified columns or expressions, which can significantly enhance query efficiency [10].

In the context of funnel analysis, clustering the Events table by event\_time can yield several benefits:

- Reduced I/O: Queries filtering on event\_time can skip reading irrelevant data blocks
- Improved sort performance: Data may already be partially or fully sorted, reducing the computational cost of window function operations
- Enhanced data locality: Related events are stored together, potentially improving cache hit rates

To implement effective clustering for funnel analysis:

- Choose clustering keys that align with common query patterns (e.g., event\_time and user\_id)
- Regularly maintain clustering through automated processes or manual reorganization
- Monitor query performance to ensure clustering continues to provide benefits as data and query patterns evolve

While clustering can offer significant performance improvements, it may not be equally beneficial for all queries. Queries that don't align with the clustering strategy may not see the same level of performance enhancement.

When deciding between these approaches and considering their respective optimizations, it's essential to:

- 1. Analyze your specific workload and query patterns
- 2. Conduct performance testing with representative data volumes
- 3. Consider the trade-offs between query performance, storage costs, and maintenance overhead
- 4. Regularly reassess your approach as data volumes grow and query patterns evolve

By carefully considering these performance optimizations and tailoring them to your specific use case, you can ensure that your funnel analysis queries execute efficiently, providing timely insights to drive product decisions.

Characteristic	Join Sequence with Pre- computed Indices	Stacked Window Functions with Table Clustering
Query Performance	High	High
Storage Overhead	High	Medium
Maintenance Complexity	High	Medium
Data Ingestion Impact	Medium	Low
Flexibility for Various Queries	Medium	High

**Table 2:** Performance Trade-offs in Funnel Analysis Approaches [9, 10]

## **CONCLUSION**

In conclusion, funnel analysis emerges as a pivotal tool for product managers, offering invaluable insights into user behavior and product performance. Implementing funnel analysis in modern data warehouses, through Join Sequence or Stacked Window Functions approaches provides a powerful means to track and optimize user journeys. Organizations can extract maximum value from their funnel analysis efforts by leveraging advanced query optimization techniques and carefully considering performance implications.

As data analytics evolves, staying informed about the latest optimization strategies and regularly reassessing analytical approaches will be crucial for maintaining a competitive edge. Ultimately, the effective use of funnel analysis, supported by robust data warehouse capabilities, empowers businesses to make data-driven decisions that enhance user experience, increase conversions, and drive overall product success.

#### REFERENCES

- [1] G. S. Day, "The capabilities of market-driven organizations," Journal of Marketing, vol. 58, no. 4, pp. 37-52, 1994. Available: https://doi.org/10.1177/002224299405800404
- [2] A. J. Alvarez, "Data-Driven Product Management: How to Use Data to Develop, Launch and Grow Your Products," O'Reilly Media, Inc., 2022. Available: https://www.oreilly.com/library/view/data-driven-product-management/9781098141325/
- [3] D. T. Bourgeois and T. Bourgeois, "Information Systems for Business and Beyond," Pressbooks, 2019. Available: https://bus206.pressbooks.com/chapter/chapter-2-information-systems-for-business-and-beyond/
- [4] A. Kohavi, A. Tang, and Y. Xu, "Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing," Cambridge University Press, 2020. Available: https://experimentguide.com/
- [5] M. Kleppmann, "Designing Data-Intensive Applications: The Big Ideas Behind Reliable, Scalable, and Maintainable Systems," O'Reilly Media, 2017. Available: https://www.oreilly.com/library/view/designing-data-intensive-applications/9781491903063/
- [6] J. M. Hellerstein, M. Stonebraker, and J. Hamilton, "Architecture of a Database System," Foundations and Trends in Databases, vol. 1, no. 2, pp. 141-259, 2007. Available: https://dl.acm.org/doi/10.1561/1900000002
- [7] S. Chaudhuri, "An Overview of Query Optimization in Relational Systems," in Proceedings of the Seventeenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, 1998, pp. 34-43. Available: https://dl.acm.org/doi/10.1145/275487.275492
- [8] A. Thusoo et al., "Hive: A Warehousing Solution Over a Map-Reduce Framework," Proceedings of the VLDB Endowment, vol. 2, no. 2, pp. 1626-1629, 2009. Available: https://dl.acm.org/doi/10.14778/1687553.1687609

## Optimizing Funnel Analysis in Modern Data Warehouses

- [9] P. O'Neil and D. Quass, "Improved Query Performance with Variant Indexes," in Proceedings of the 1997 ACM SIGMOD International Conference on Management of Data, 1997, pp. 38-49. Available: https://dl.acm.org/doi/10.1145/253260.253268
- [10] D. Abadi, P. Boncz, S. Harizopoulos, S. Idreos and S. Madden, "The Design and Implementation of Modern Column-Oriented Database Systems," Foundations and Trends in Databases, vol. 5, no. 3, pp. 197-280, 2013. Available: https://doi.org/10.1561/1900000024

**Citation:** Satyam Shekhar, Optimizing Funnel Analysis in Modern Data Warehouses, International Journal of Computer Engineering and Technology (IJCET), 15(4), 2024, pp. 516-526

**Abstract Link:** https://iaeme.com/Home/article\_id/IJCET\_15\_04\_045

#### **Article Link:**

https://iaeme.com/MasterAdmin/Journal\_uploads/IJCET/VOLUME\_15\_ISSUE\_4/IJCET\_15\_04\_045.pdf

**Copyright:** © **2024** Authors. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

This work is licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0).



**☑** editor@iaeme.com