# 基数估计文献综述

张亚楠 22451207

**Abstract:**

Cardinality estimation is a critical element of query optimization in database management systems (DBMSs). Accurate predictions of tuple counts for SQL queries are essential for devising efficient execution plans that improve overall query performance. This review explores recent advancements in cardinality estimation, focusing on the Attention-based Learned Cardinality Estimator (ALECE) and its use of attention mechanisms. Additionally, other deep learning-based approaches, including NeuroCard, DeepDB, and FLAT, are discussed, highlighting their unique methodologies and contributions to this evolving field.

**Introduction:**

Traditional methods for cardinality estimation, such as histogram-based and sampling-based approaches, have been limited in their ability to capture complex data patterns and attribute dependencies. With the rise of machine learning (ML) and deep learning (DL), new methods have been developed to address these limitations. This review focuses on four recent deep learning-based models—ALECE, NeuroCard, DeepDB, and FLAT—that exemplify the latest advancements in cardinality estimation for SQL queries, each providing novel solutions to the challenges of accuracy and efficiency in database query optimization.

## 1.ALECE: Attention-based Learned Cardinality Estimator

ALECE, or the Attention-based Learned Cardinality Estimator, distinguishes itself as a highly effective model for SQL query cardinality estimation by applying sophisticated attention mechanisms that capture nuanced dependencies between query patterns and underlying database attributes. This approach enables ALECE to excel where traditional estimation methods often fall short, particularly with complex or highly interdependent data.

**Data-Encoder Module:** The data-encoder module in ALECE is specifically designed to handle intricate attribute relationships within the data distribution by using self-attention layers. These layers enable ALECE to learn correlations among attributes, which are often missed by traditional cardinality estimation models. By capturing these interdependencies, the data-encoder produces a contextually enriched representation of the database that reflects the joint distributions across all attributes. This foundation strengthens the model's ability to interpret the structural

and relational intricacies of the data, setting a robust stage for accurate query analysis.

**Query-Analyzer Module:** ALECE's query-analyzer module then takes this enriched data encoding and combines it with query-specific features to produce a query-aware representation. It employs attention mechanisms to selectively focus on the data elements most relevant to the query, effectively filtering out noise and prioritizing pertinent data segments. The output of this module is a fixed-dimensional "answering" vector, which captures the essential information required to produce an accurate cardinality estimate. This vector is then processed through a simple linear regression layer that translates the encoded data-query relationship into a final cardinality estimate. By aligning data and query representations through attention, ALECE adapts seamlessly to diverse data distributions and complex query patterns, making it well-suited for dynamic and high-dimensional datasets.

**Featurization of Data and Queries:** A key strength of ALECE lies in its meticulous approach to featurizing both data and queries, allowing it to better align the representations with real-world database characteristics. Data attributes are transformed into normalized histogram vectors that fall within a [0, 1] range, enabling consistent handling across different scales and distributions. Temporal aspects or other contextual features can also be incorporated, broadening the model's adaptability. Queries are represented by boundary vectors for each filter predicate, also normalized, effectively defining a "query hyperrectangle" within the data space. This systematic design allows ALECE to model not only the data distribution but also the scope of the query with high precision, providing a robust framework for nuanced cardinality estimation.

**Impact and Benchmarking:** Through these innovations, ALECE achieves a high degree of accuracy, significantly outperforming traditional models and setting a new benchmark in the field of learned cardinality estimation. The combination of self-attention within the data encoder and selective attention in the query analyzer enables ALECE to address the complex dependencies and patterns present in large databases, especially for multi-table and high-dimensional queries. This model represents a new frontier in database query optimization, exemplifying how attention mechanisms can be harnessed to capture the deep structure of both data and query representations.

## 2. NeuroCard

NeuroCard introduces an efficient approach to cardinality estimation by focusing on innovative sampling and factorization techniques, reducing computational demands while preserving accuracy. This model is particularly valuable in large-scale database environments, where traditional cardinality estimation techniques are often impractical due to the sheer volume of data and complex join structures.

NeuroCard's design not only improves efficiency but also enhances scalability, making it well-suited for high-dimensional, high-cardinality datasets.

**Join Sampling:** A core feature of NeuroCard is its ability to approximate join operations through selective sampling rather than exhaustive computation. By carefully sampling joins rather than fully calculating them, NeuroCard drastically reduces training time and computational costs, which are often major bottlenecks in large databases. This sampling technique balances efficiency and accuracy, allowing the model to generalize well across different query types without requiring the costly computation of all potential join combinations.

**Single Model for Multiple Queries:** Unlike models that necessitate separate training for distinct queries or query types, NeuroCard is designed to handle multiple queries with a single, unified model. This approach provides significant flexibility by enabling NeuroCard to process any subset of tables or joins within the database without the need for re-training or fine-tuning for each specific query. This single-model architecture reduces maintenance efforts and computational costs, making NeuroCard a practical choice for dynamic database environments where query patterns may change frequently.

**Lossless Column Factorization:** Handling high-cardinality columns is a common challenge in cardinality estimation, as they can inflate model size and computational demands. NeuroCard addresses this with a lossless column factorization method, which represents these columns efficiently without losing information about their relationships to other attributes. This factorization approach captures the essential dependencies within high-cardinality columns while controlling model size, ensuring that the model scales well even with large and complex datasets. By preserving data relationships, NeuroCard maintains the integrity of the database schema, enabling it to generalize effectively across diverse data structures and query types.

**Scalability and Efficiency:** Through its combination of join sampling, flexible single-model architecture, and efficient column factorization, NeuroCard achieves a high degree of scalability. This design allows it to manage large and complex databases efficiently, making it a cost-effective solution for scenarios where full join computation is impractical or prohibitively expensive. NeuroCard's adaptability to high-dimensional data structures and its robust performance in multi-table query contexts further underscore its value in database query optimization.

### 3. DeepDB

DeepDB is a versatile deep learning framework for cardinality estimation that excels in adapting to a wide variety of query types and data configurations. Unlike traditional models that rely heavily on predefined statistical structures, DeepDB is specifically designed to handle the dynamic, heterogeneous nature of modern databases. Through its flexible architecture and efficient handling of large-scale

datasets, DeepDB provides a practical and adaptable solution for environments where query patterns are highly variable and data structures are complex.

**Model Flexibility:** A key strength of DeepDB lies in its ability to generalize across diverse query types and database schemas. This flexibility is achieved by decoupling the model from rigid structures, allowing it to adapt dynamically to different database configurations. DeepDB can handle complex query patterns that span multiple tables and involve various joins and aggregations, making it suitable for highly relational data scenarios. Its adaptability enables it to accommodate evolving schemas and changes in data distributions without extensive retraining, giving it a significant advantage over more traditional, rigid estimation methods.

**Efficiency in Large Databases:** DeepDB is engineered to operate efficiently within large-scale database systems, where both data volume and query frequency may be substantial. Its architecture is optimized to handle frequent query updates and to manage data with high dimensionality or a wide range of attribute types. By leveraging advanced learning techniques, DeepDB minimizes the computational resources needed for training and inference, making it particularly well-suited for environments that demand real-time or near-real-time response rates. This optimization enables DeepDB to deliver accurate cardinality estimates without sacrificing performance, even as database sizes grow or when handling high-throughput workloads.

**Adaptability and Practical Advantages:** DeepDB's adaptability allows it to outperform traditional histogram- and sampling-based approaches, particularly in scenarios where query patterns are unpredictable or rapidly changing. By bypassing the limitations of static distribution assumptions, DeepDB can accurately capture complex data relationships that other methods might overlook. Its use of deep learning allows it to model non-linear dependencies and intricate attribute interactions more effectively, providing more reliable estimates in varied and dynamic database environments. This makes DeepDB a valuable choice for applications where flexibility, scalability, and consistent accuracy are essential.

### 4. FLAT (Factorize-Split-Sum-Product Network)

FLAT, or the Factorize-Split-Sum-Product Network, introduces a unique and effective approach to cardinality estimation by breaking down the estimation task into smaller, more manageable sub-problems. This factorized approach not only streamlines the computation process but also enhances the model's accuracy and adaptability when handling complex, multi-attribute queries in high-dimensional data environments. FLAT's architecture is specifically designed to balance computational efficiency with precision, making it an attractive choice for databases with intricate join structures and diverse data relationships.

**Factorization through FSPN:** At the core of FLAT's design is the Factorize-Split-Sum-Product Network (FSPN) architecture, which decomposes the cardinality estimation task into distinct sub-tasks. By factorizing the problem into smaller components, FLAT transforms a potentially unwieldy cardinality estimation process into a series of tractable computations. These sub-problems are processed individually and later aggregated through a summing and product mechanism, combining the sub-estimates to produce a final, cohesive cardinality estimate. This modular approach allows FLAT to efficiently capture complex dependencies within the data, while maintaining computational feasibility and minimizing the risk of estimation errors.

**Handling Complex Queries Efficiently:** FLAT's factorized architecture is particularly advantageous for complex queries involving multiple joins, intricate filter predicates, and high-cardinality columns. By decomposing each query into its component factors, FLAT can effectively manage the complexity without becoming computationally intensive. This factorization method enables the model to process high-dimensional queries with relative ease, as each sub-component of the query can be handled independently before combining them into a final estimate. This approach allows FLAT to remain responsive and accurate even with challenging, large-scale queries, making it well-suited for dynamic and data-rich environments where query structures are often intricate and multi-layered.

**Achieving Balance Between Accuracy and Computational Efficiency:** Through its innovative factorization approach, FLAT achieves a balanced trade-off between accuracy and computational efficiency. Unlike traditional methods, which can struggle with high-dimensional data or require extensive resources for complex queries, FLAT's FSPN architecture provides a scalable solution that remains both accurate and efficient. The model's ability to factorize complex dependencies without incurring excessive computational costs makes it a practical choice for databases with diverse, high-dimensional data. This design ensures that FLAT can deliver precise cardinality estimates consistently, while maintaining manageable levels of computation, even as data volume or query complexity grows.

### 5. Sample-Efficient Cardinality Estimation Using Geometric Deep Learning

Sample-efficient cardinality estimation leverages geometric deep learning to deliver accurate query estimates, even in scenarios with limited data. These methods capitalize on graph neural networks (GNNs) and principles of geometric deep learning to understand the complex relational structures of high-dimensional database schemas while significantly reducing dependency on large data samples. This approach is particularly advantageous in environments with sparse or incomplete data distributions, where traditional methods may struggle to maintain accuracy. Key aspects of this approach include:

**Graph Representation of Database Relations:** Sample-efficient geometric models treat databases as graph structures, where tables, attributes, and even relationships are represented as nodes and edges. This graphical representation mirrors the relational nature of database schemas, allowing the model to efficiently encode structural and relational information. By capturing the inherent multi-relational dependencies within the data, this method provides a robust framework for learning the relational context needed for accurate cardinality estimation. Graph-based representation thus enables the model to address complex query structures by naturally incorporating dependencies among different database entities.

**Message Passing for Relational Context Understanding:** Similar to ALECE's Join Graph Message Passing, geometric deep learning methods employ message-passing techniques within graph neural networks to propagate information across nodes. This allows each node (representing an attribute, table, or relationship) to gather information from its neighbors and form a contextually enriched representation of the relational structure. However, geometric deep learning models extend this capability by enhancing the model's generalization across diverse and complex schemas, often with fewer samples than traditional methods would require. This ability to generalize effectively, even across complex and high-dimensional relationships, ensures that the model can capture the nuances of relational data without excessive computational overhead.

**Reduced Sample Dependency:** One of the most compelling advantages of geometric deep learning in cardinality estimation is its sample efficiency. By modeling inherent data relationships through graph structures, these models minimize reliance on large datasets, making them well-suited for cases where data acquisition is costly or samples are limited. The use of graph representations allows the model to extrapolate patterns from smaller amounts of data, enhancing its robustness and accuracy in resource-constrained scenarios. This reduction in sample dependency is particularly beneficial for databases with sparse distributions, where obtaining representative samples for training can be challenging.

**Advantages in Resource-Constrained Environments:** Sample-efficient geometric deep learning models offer a promising alternative for environments with limited computational or data resources. By combining high accuracy with minimal sample requirements, these models can maintain performance levels comparable to data-intensive methods while operating within the constraints of limited datasets. This makes them highly valuable for applications in resource-constrained environments, where traditional approaches may struggle due to their reliance on large training datasets and high computational demands.

**Conclusion:**

The landscape of cardinality estimation has evolved significantly with the advent of deep learning-based methods. Models such as ALECE, NeuroCard, DeepDB, and FLAT

represent cutting-edge approaches that enhance accuracy and efficiency in query optimization. Each model brings unique solutions to the challenges posed by diverse database structures and query patterns, setting a foundation for further advancements in DBMS performance and reliability. As this field progresses, these deep learning-driven techniques hold immense potential for transforming the way queries are optimized in modern database management systems.