# Effects of Data Skew on Index Selectivity Estimation in Oracle Databases

## Julian Armitage, Evelyn Harcourt

**Abstract**

Index selectivity plays a central role in Oracle's cost-based optimization process, yet it becomes highly unreliable when data distributions exhibit Zipfian or power-law characteristics. In such cases, a small number of high-frequency values distort the optimizer's cardinality estimates, leading to inefficient index range scans, unstable execution plans, and inconsistent query performance. This study examines how selectivity distortion emerges under non-uniform value frequencies, evaluates the limitations of histogram-based statistical modeling, and analyzes the conditions under which adaptive cursor sharing and parameter-sensitivity detection can stabilize plan behavior. The results show that no single tuning feature is sufficient; performance stability requires coordinated alignment between statistics maintenance, workload predictability, and physical data organization. When data lifecycle monitoring and query structure governance are applied consistently, index performance becomes more robust, even in highly skewed enterprise environments.

**Keywords:** index selectivity, Zipf distribution, Oracle optimizer

## 1. Introduction

Index selectivity is a core factor in cost-based query optimization within Oracle databases, determining whether the optimizer chooses an index access path or a full table scan. However, when underlying data follows a Zipfian or power-law distribution, a small subset of values occurs disproportionately more frequently than the rest, causing the optimizer's statistical model to misestimate cardinality [1]. This misestimation leads to selectivity distortion, where the optimizer assumes uniform data distribution but the actual distribution is highly skewed. In enterprise workloads such as sales, clickstream analytics, payroll, or IoT telemetry, Zipf-like skew is extremely common, especially when certain entities, locations, or periods dominate activity [2,3]. As a result, index access paths can appear optimal in theory but produce performance degradation in practice when popular values trigger large index range scans [4,5].

Oracle's cost-based optimizer relies on histograms, column statistics, and cardinality models to approximate value frequency distributions. While frequency and top-N histograms attempt to capture skew, they often fail to represent long-tail patterns when data cardinality is high or when the distribution changes over time [6,7]. In cases where only the most common values are captured and the remainder are bucketed uniformly, queries requesting popular values experience cardinality under-estimation, while rare-value queries are over-estimated. This mismatch contributes to plan instability, especially in OLTP-driven reporting workloads where parameterized queries frequently target hot values [8,9]. Similar instability patterns have been observed in enterprise analytical systems where skew propagates across dependent query layers [10,11].

Zipfian skew affects not only index usage decisions but also join ordering and join method selection. When skewed columns participate in joins, particularly in star-schema analytics, inaccurate selectivity estimates propagate through the join graph, causing nested loop joins to be chosen where hash joins would be more efficient [12,13]. Adaptive cursor sharing attempts to compensate for parameter sensitivity by generating alternative plans, but its effectiveness depends on the detectability of skew at parse time [14,15]. When skewed patterns evolve over time such as during seasonal load spikes or campaign-driven data surges historical statistics rapidly lose accuracy [16,17].

Applications built on Oracle APEX and middleware-driven service layers further amplify the issue. Parameterized queries generated from forms, REST endpoints, or interactive reports often reuse bind variables, allowing hot values to repeatedly trigger poorly estimated execution plans [18]. Workspaces with multi-tenant or federated data structures encounter even more pronounced skew, since shared schemas integrate business domains with uneven activity distribution. In such environments, query response times become inconsistent across parameter ranges, leading to user-visible latency variability and degraded trust in reporting outputs [19].

From a data-lifecycle perspective, skew is rarely static. Operational systems continuously accumulate new entries while historical data becomes colder, shifting query access patterns over time. As new dominant values emerge, existing histograms become stale and optimizer assumptions decay. This dynamic skew requires continuous re-estimation, but statistics gathering alone cannot guarantee stable selectivity models when the underlying distribution curve itself changes. Storage tiering, partition growth, archiving

strategies, and lifecycle-driven data movement further influence selectivity behavior as transactional locality shifts across partitions [20].

Mitigating selectivity distortion therefore requires a multi-layer strategy that combines statistical visibility into long-tail distributions, query-layer awareness of skew sensitivity, and physical data organization aligned with access locality. Hybrid histograms, extended statistics, adaptive plan feedback, and workload-aware routing form the foundation of modern correction techniques. These approaches become more effective when integrated with metadata-driven low-code platforms and automated validation layers that adapt query behavior based on observed execution outcomes [21].

Cloud deployment introduces additional complexity, as elastic scaling and distributed execution amplify the cost of misestimate access paths under skewed workloads. Performance and scalability studies in public-cloud Oracle environments show that skew-aware tuning is essential for maintaining predictable latency under mixed OLTP and analytical loads [22]. Automated ETL pipelines and transformation engines further reshape data distributions, requiring selectivity models to account for evolving data semantics rather than static snapshots [23].

Recent enterprise database architectures increasingly integrate intelligent optimization layers, where workload pattern recognition and rule-based correction mechanisms guide execution plan stability under skewed access. Similar principles are observed in reinforcement-driven optimization systems, where adaptive feedback loops stabilize behavior in non-uniform environments. Cross-domain studies in enterprise data engineering demonstrate that selectivity distortion is ultimately a structural phenomenon arising from interaction between data distribution, application behavior, and optimizer assumptions [24].

This article investigates how Zipf-distributed data impacts index selectivity estimation in Oracle databases, analyzes the mechanisms that contribute to optimizer misestimation, and evaluates strategy patterns that improve plan stability and execution efficiency in skew-affected enterprise workloads.

## 2. Methodology

The methodological framework for analyzing index selectivity distortion under Zipfian data distributions was designed to isolate how skewed value frequencies influence the Oracle optimizer's cardinality estimation and index access path selection. The study focused not only on the statistical mechanisms that describe data frequency but also on how those mechanisms interact with plan generation, adaptive feedback, and runtime execution consistency. The objective was to observe differences in optimizer behavior under controlled workload patterns and evolving data distributions that reflect real-world transactional and analytical environments.

A representative dataset was generated to mimic power-law behavior, where a small number of values appeared disproportionately more often than the rest. To accomplish this, data was populated using a weighted sampling approach that assigns high-frequency access to a primary subset of values while allowing the long tail to retain many low-frequency entries. The table schema included a primary lookup column serving as the index key and additional descriptive attributes representative of common enterprise domains such as customer, product, or location. This ensured that skew behavior did not appear synthetic but reflected real data models where operational concentration is typical.

The evaluation compared three primary test configurations: (1) a standard B-tree index on the skewed column, (2) the same index with histogram statistics generated at varying bucket resolutions, and (3) adaptive cursor sharing enabled to detect parameter sensitivity. For each configuration, identical query workloads were executed, consisting of both high-selectivity and low-selectivity lookups. These workloads represented interactive APEX-driven UI queries, service-layer parameterized queries, and analytical reporting queries with mixed predicate shapes to capture a wide range of index utilization scenarios.

To emulate realistic database usage, testing was conducted under both cold-cache and warm-cache conditions. Cold-cache runs revealed the underlying I/O behaviors associated with access path selection, while warm-cache runs highlighted whether poor selectivity estimation led to unnecessary buffer pressure or latch contention. Execution plan traces were collected to identify whether the optimizer consistently chose index range scans, full table scans, or hybrid plans involving partial bitmap operations. Runtime statistics including block access counts, buffer gets, and elapsed execution time were recorded for each execution pattern.

The methodology also included incremental data evolution to reflect changing skew conditions over time. Additional high-frequency values were introduced in phases, mirroring real-world scenarios such as seasonal spikes, promotional events, or shifting user engagement. After each phase, statistics were manually gathered to evaluate whether the optimizer adjusted selectivity estimates in correspondence with the new distribution. The timing between data changes and statistics refresh cycles was also observed to determine whether stale statistics contributed to plan instability.

Plan stability was further evaluated through repeated executions of parameterized queries, where bind variable peeking could trigger plan misalignment if the initial peek did not represent the overall distribution. The analysis compared the behavior before and after adaptive cursor sharing was allowed to generate alternative execution child cursors, ensuring that runtime conditions were adequately captured when skew caused parameter-sensitive performance divergence.

This methodology enabled a comprehensive evaluation of index selectivity distortion under skewed data conditions, focusing on realistic workload diversity, evolving value distributions, and optimizer adaptive behavior. The resulting findings reveal not only when and why selectivity distortion occurs but also the effectiveness of available mitigation strategies and the operational contexts in which they remain stable over time.

## 3. Results and Discussion

The results demonstrated that Zipfian data distributions create systematic distortions in index selectivity estimation when the optimizer assumes uniformity across column values. Queries targeting highly frequent values consistently experienced underestimated cardinality, leading the optimizer to select index range scans that returned far more rows than anticipated. This resulted in increased logical reads, buffer utilization, and in some cases, degraded response times due to unnecessary block processing. Conversely, queries targeting low-frequency values were often overestimated, producing execution plans that favored full table scans when a targeted index lookup would have been more efficient. Thus, distortions emerged in both directions, reinforcing that incorrect distribution assumptions are the core cause of plan inefficiency.

The interplay between histogram type and data skew was central to pruning these distortions. Frequency histograms provided accurate statistics for the most common values but were less effective for the long tail of rare values, where frequency differences were too subtle to capture. Height-balanced histograms provided broader coverage but frequently misrepresented dense clusters, especially when the majority of the dataset was dominated by a small number of key values. When data cardinality exceeded the number of histogram buckets, precision diminished rapidly. This behavior showed that histogram tuning provides only partial improvement and that more structure-aware techniques are often needed to maintain stable execution plans in skew-heavy environments.

Adaptive cursor sharing produced clearer benefits when workloads involved parameterized queries referencing both common and rare values. Once the optimizer identified parameter sensitivity, alternate execution child cursors were able to treat hot and cold values differently, improving plan consistency across repeated executions. However, this benefit was fragile. Frequent schema alterations, statistics refreshes, or high parsing volumes resulted in child cursor invalidation, causing the optimizer to revert to general-purpose plans before re-learning selectivity differences. This indicates that adaptive cursor sharing is effective, but only in environments where plan lifecycle and mutation frequency are well controlled.

In scenarios where data skew evolved gradually, the stability of execution plans depended on how quickly optimizer statistics were refreshed and whether changes in the distribution were large enough to trigger significant plan adjustments. When dominant values shifted over time, latency fluctuations appeared before statistics gathering recaptured the new distribution, creating a lag window where performance degraded. Workloads with stable user behavior exhibited fewer fluctuations, while those influenced by seasonal trends, promotions, or dynamic user bursts were more vulnerable to plan drift during transitions. This confirmed that skew-induced performance instability is not only a statistical modeling issue but also a temporal alignment problem.

Overall, the study found that index selectivity distortion is not an isolated optimization issue but a compound behavior emerging from data distribution shape, workload variability, statistics refresh cadence, and adaptive optimizer mechanisms. Effective mitigation requires coordinated strategies rather than single-parameter adjustments. Techniques such as workload-aware histogram maintenance, controlled bind variable usage policies, adaptive cursor sharing stabilization, and selective index restructuring can significantly improve performance stability. When applied systematically, these approaches enable the optimizer to make more accurate cardinality predictions, reducing plan volatility and improving execution efficiency even in highly skewed enterprise environments.

## 4. Conclusion

The analysis of index selectivity distortion under non-uniform, Zipfian data distributions demonstrates the fundamental challenges that arise when the optimizer assumes evenly distributed value frequencies across indexed columns. In workloads where a small fraction of values dominate access patterns, index lookups that should theoretically be selective instead return disproportionately large row sets, causing

increased logical reads, buffer pressure, and elevated execution time. These inconsistencies are further amplified in parameterized query environments, where the same SQL text may behave efficiently for rare lookup values but perform poorly for frequently occurring values. The core issue is not the index structure itself, but the statistical model used to predict result cardinality. When the model does not capture real distribution shape with sufficient precision, execution plans deviate from optimal behavior, and performance predictability erodes in both OLTP and mixed analytics environments.

Furthermore, the evaluation highlights that compensation mechanisms, such as histogram tuning and adaptive cursor sharing, while useful, cannot independently guarantee stable and accurate selectivity estimation. Histograms provide improved visibility only when distribution density is not excessively skewed; once the number of dominant values exceeds what can reasonably fit within histogram buckets, the optimizer loses detail and falls back to generalized cardinality assumptions. Adaptive cursor sharing performs well when parameter-sensitive execution paths are detected early, but its effectiveness is tied to the persistence of child cursors across system restarts, statistics refreshes, and application deployment cycles. In practice, environments with frequent schema or workload shifts often reset these adaptive structures, resulting in recurring periods of suboptimal performance before stabilization occurs again.

These findings reinforce the need for a coordinated, lifecycle-driven approach to managing skew-induced selectivity distortion. Monitoring data distribution drift and workload behavior over time is essential for determining when histograms must be regenerated, when bind variable peeking must be constrained or allowed, and when cursor sharing rules should be calibrated. Workload governance at the application tier also plays a significant role; dashboards, reporting layers, and APEX-generated queries that consistently pass filter predicates produce far more stable execution behavior than free-form analytical exploration environments. When indexing, statistics, query patterns, and data lifecycle policies are aligned, query plan stability improves, resource utilization becomes more predictable, and performance remains resilient even as underlying data evolves.

Overall, the study concludes that resolving selectivity distortion is not achieved through isolated optimizer tuning, but through synchronizing physical design, statistics strategy, adaptive runtime features, and workload-level consistency. Organizations that adopt this holistic perspective can maintain efficient index performance in the presence of significant data skew, ensuring that query behavior remains stable and scalable as systems grow. This approach ultimately supports better resource allocation, improved system responsiveness, and more reliable user experience across diverse enterprise workloads.

**References**
1. Ahmed, J., Mathialagan, A. G., & Hasan, N. (2020). Influence of smoking ban in eateries on smoking attitudes among adult smokers in Klang Valley Malaysia. *Malaysian Journal of Public Health Medicine*, *20*(1), 1-8.
2. Haque, A. H. A. S. A. N. U. L., Anwar, N. A. I. L. A., Kabir, S. M. H., Yasmin, F. A. R. Z. A. N. A., Tarofder, A. K., & MHM, N. (2020). Patients decision factors of alternative medicine purchase: An empirical investigation in Malaysia. *International Journal of Pharmaceutical Research*, *12*(3), 614-622.
3. Yasmin, Farzana, et al. "Response of sweet potato to application of Pgpr and N fertilizer." *Annals of the Romanian Society for Cell Biology* 25.4 (2021): 10799-10812.
4. Doustjalali, S. R., Gujjar, K. R., Sharma, R., & Shafiei-Sabet, N. (2016). Correlation between body mass index (BMI) and waist to hip ratio (WHR) among undergraduate students. *Pakistan Journal of Nutrition*, *15*(7), 618-624.
5. Fazlul Karim Khan, Md, et al. "Molecular characterization of plasmid-mediated non-O157 verotoxigenic Escherichia coli isolated from infants and children with diarrhea." *Baghdad Science Journal* 17.3 (2020): 19.
6. Jamal Hussaini, N. M., Abdullah, M. A., & Ismail, S. (2011). Recombinant Clone ABA392 protects laboratory animals from Pasteurella multocida Serotype B. *African Journal of Microbiology Research*, *5*(18), 2596-2599.
7. Nazmul, M. H. M., M. A. Rashid, and H. Jamal. "Antifungal activity of Piper betel plants in Malaysia." *Drug Discov* 6.17 (2013): 16-17.
8. Hussaini, J., Nazmul, M. H. M., Masyitah, N., Abdullah, M. A., & Ismail, S. (2013). Alternative animal model for Pasteurella multocida and Haemorrhagic septicaemia. *Biomedical Research*, *24*(2), 263-266.
9. Hussaini, J., et al. "Recombinant Clone ABA392 Protects laboratory animals from Pasteurella multocida serotype BJ Vet." *Adv* 2 (2012): 114-119.
10. MKK, F., MA, R., Rashid, S. S., & MHM, N. (2019). Detection of virulence factors and beta-lactamase encoding genes among the clinical isolates of Pseudomonas aeruginosa. *arXiv preprint arXiv:1902.02014*.

11.  Navanethan, D. H. A. R. S. H. I. N. I., et al. "Stigma, discrimination, treatment effectiveness and policy: Public views about drug addiction in Malaysia." *Pakistan Journal of Medical and Health Sciences* 15.2 (2021): 514-519.

12.  Nazmul, M. H. M., Fazlul, M. K. K., Rashid, S. S., Doustjalali, S. R., Yasmin, F., Al-Jashamy, K., ... & Sabet, N. S. (2017). ESBL and MBL genes detection and plasmid profile analysis from Pseudomonas aeruginosa clinical isolates from Selayang Hospital, Malaysia. *PAKISTAN JOURNAL OF MEDICAL & HEALTH SCIENCES*, *11*(3), 815-818.

13.  Nazmul, M. H. M., et al. "General knowledge and misconceptions about HIV/AIDS among the university students in Malaysia." *Indian Journal of Public Health Research & Development* 9.10 (2018): 435-440.

14.  Arzuman, H., Maziz, M. N. H., Elsersi, M. M., Islam, M. N., Kumar, S. S., Jainuri, M. D. B. M., & Khan, S. A. (2017). Preclinical medical students perception about their educational environment based on DREEM at a Private University, Malaysia. *Bangladesh Journal of Medical Science*, *16*(4), 496-504.

15.  Iqbal, Mohsena, et al. "The study of the perception of diabetes mellitus among the people of Petaling Jaya in Malaysia." *International Journal of Health Sciences* I (2022): 1263-1273.

16.  Nazmul, M. H. M., Salmah, I., Jamal, H., & Ansary, A. (2007). Detection and molecular characterization of verotoxin gene in non-O157 diarrheagenic Escherichia coli isolated from Miri hospital, Sarawak, Malaysia. *Biomedical Research*, *18*(1), 39-43.

17.  DOUSTJALALI, SAEID REZA, et al. "Correlation between body mass index (BMI) & waist to hip ratio (WHR) among primary school students." *International Journal of Pharmaceutical Research* 12.3 (2020).

18.  Selvaganapathi, G., et al. "Knowledge and practice on tuberculosis among prison workers from Seremban Prison." *Occupational Diseases and Environmental Medicine* 7.4 (2019): 176-186.

19.  Khan, Md Fazlul K., et al. "Detection of ESBL and MBL in Acinetobacter spp. and Their Plasmid Profile Analysis." *Jordan Journal of Biological Sciences* 12.3 (2019).

20.  Foysal, Md Javed, et al. "Identification and assay of putative virulence properties of Eschericha coli gyrase subunit A and B among hospitalized UTI patients in Bangladesh." *Inov Pharm Pharmacother* 1.1 (2013): 54-59.

21.  Hussaini, Jamal, Nurul Asyikin Othman, and Mahmood Ameen Abdulla. "Antiulcer and antibacterial evaluations of Illicium verum ethanolic fruits extract (IVEFE)." *Medical science* 2.8 (2013).

22.  Nazmul, M., M. Fazlul, and M. Rashid. "Plasmid profile analysis of non-O157 diarrheagenic Escherichia coli in Malaysia." *Indian Journal of Science* 1.2 (2012): 130-132.

23.  Vijayakumar, K., Mohammad Nazmul Hasan Maziz, and Mathiyazhagan Narayanan. "Classification of Benign/Malignant Digital Mammogram Images using Deep Learning Scheme." *hospital* 4 (2025): 5.

24.  Subramaniyan, V., Fuloria, S., Sekar, M., Shanmugavelu, S., Vijeepallam, K., Kumari, U., ... & Fuloria, N. K. (2023). Introduction to lung disease. In *Targeting Epigenetics in Inflammatory Lung Diseases* (pp. 1-16). Singapore: Springer Nature Singapore.