

Learned Indexes

Machine learning in databases:

07.03.18

- HashTables
- B-Trees

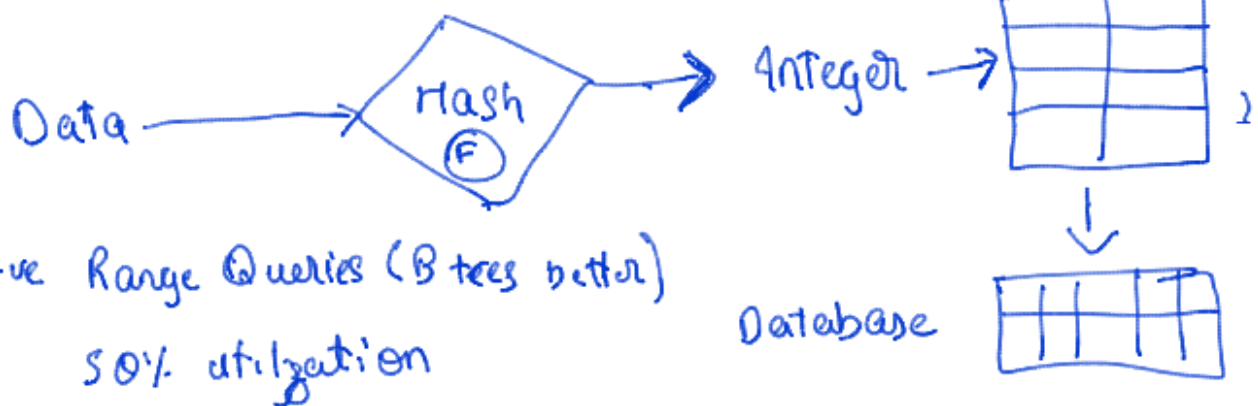
Tyler Elnot Betilyon

Performance of Indexes

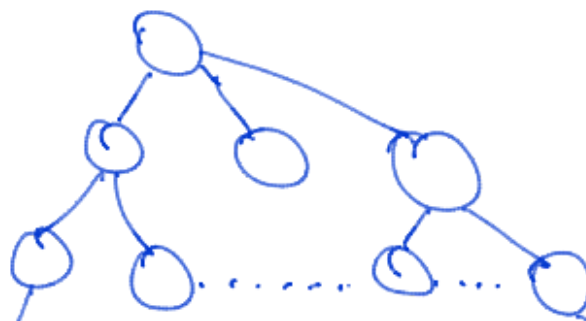
1. Type of Query (Range - exact)
2. Lookup Speed
3. Cost to delete
4. Size & Utilization
5. Minimise (redo) errors rather than Operations

HashTable

- Collision handling technique
Cuckoo hashing



B-Tree



○ First

○ Last

- Sorted
 - Range Queries
 - Query time $(\log n)$ / Search & Delete
 - Must be balanced
 - Split Operation (Insertion, Deletion) - Extra Cost
-

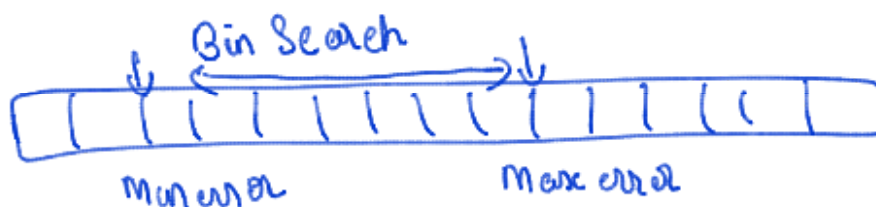
Learned Indexes

- Estimation \Rightarrow (func) * Key \Rightarrow Regression
- ML can be used?



Adv.

- Takes into account distribution of data



-
- Small fully connected NN.

- 4n m/m⁰

Btree vs Learned Indexes

- Speedup
- Better space improvement
- Save node size

Results & Conclusion

models take into account distribution

- Overfitting is the goal - Accuracy
- Data Warehousing application?
- Index using need not grow using ML model. Standard index grows.

- memory \rightarrow X dead classical CPU [X parallelization]
- ML models \Rightarrow parallelization
- Multidimension data?

Q&A

- Subtree vs Btree
- Retrieval of insertion & records (train in 1 sec)
- Normal distribution preferred?

no \rightarrow model works for any distribution

- Could GANs can be used to generate indices? (cluster)
 - Parallelization in GPUs for databases is future?
 - How ML model saves data?
 - Hashmaps produces Uniform distribution
 - ML takes into acc data distribution, produces fewer hash collision.
 - Why GPU > CPU?
 - GPU has vector instructions
 - CPU does on 64 bit instruction
 - 32 bits < 64 bit instruction, clock cycle is slower
 - Process 16 times faster, but clk is 4 times slower.
-