

final review

wbg231

January 2023

1 RDMS

review

- relational model (table and data frames) standardizes data organization
- schema constrains simplify making sure data is formatted correctly
- tables can speed up common access patterns
- sql is good
- constraining data organization makes high level operations easier

question

- none :(

map reduce

- map takes input outputs key value pairs
- reducer takes key values pairs and aggregates them on key value
- combiners are optional but work work like reducers within map nodes means less data shuffling
- constraining the form of computation makes parallelism easier

question

- Combiners in map-reduce improve efficiency by decreasing the number of intermediate keys prior to the reduce phase.

- false combiners preserve the key structure but produce fewer values they fix key skew
- map function in map reduce must produce at least one intermediate output
- false

HDFS

- data nodes store partitions of large data block
- name nodes store a map of file names/and the mapping of blocks to data nodes (that is what data node has what block of data)
- data is read and append only. this makes concurrent access easier
- Replication factor more copies of each block up front means less computation and shuffling later on improves locality
- when the name node fails spark fails it is not acid

spark

- RDD are the main abstraction
- transformations take RDD to RDD and actions take RDD to results
- lineage graphs represent overall computation
- provides a higher level interface for distributed computing than map reduce

question

- compare a RDD vs a data frame
- a RDD is closer to a list it is really fast with filters and maps
- both have lazy execution
- a dataframe is made up of RDDs good for more complex analysis
- data frames in spark are read only
- true the values of a spark data frame can not be modified in place
- each step in an RDD lineage graph computation must complete before starting the next. false that is one of the reasons we have lineage graphs
- spark uses piellines to conect multiple map reduce programs
- false spark does not use map reduce

column oriented storage

- organizing data by columns makes things faster
- we are constraining data types so faster memory access
- data compression makes faster communication
- dremel is used for taking nested structured documents to tabular representations
- tables = columns = compression
- parquet is the default for spark
- parquet can be faster than spark

questions

- explain parquet and dremel
- The Dremel system was designed to efficiently process all attributes for subsets of records in a dataset
- false that is the opposite of what it does
- When written to HDFS, Parquet files locate different columns in different HDFS blocks.
- false parquet files divide blocks by column

dask

- like spark but different
- works with the scipy stack
- integrated into python directly so more flexible with different types of data questions
- spark is great when your data looks like dataframes
- when your data is not dataframes dask is good
- but more flexibility means less automatic optimization
- spark has bags that will hold it all

questions

- dask is entirely python, spark is slowed down a lot by going in and out of python
- dask is less polished
- dask is suited for working with python packages that already exist (also can work better on a single machine)
- spark is better for bad news like things but dask uses pandas directly
- HOW TO OPTIMIZE: in SQL, Spark and Dask, i.e indices (types and uses), partitions, when to use RDDs/Data frames. Also, when to use each of the frameworks.
- use filters as quickly as possible to cut down your data
- understand your lineage graph
- optimize partition structure to avoid wide dependencies

similarity search

- minhash, represent each set in a collection by the smallest hash outputs of its elements
- probability that sets collide is the jaccard similarity
- use multiple hashes to estimate jaccard similarity
- instead of partitioning use a hash function that is imperfect if there are collisions those are candidate pairs do more similarity search on those
- combine multiple min hash outputs together as a block
- block size + number of blocks can be used to boost likelihood of collision
- idea generalizes to other distance metrics

questions

- what is multi probe lsh
- multi probe lsh is instead of looking at exact matches for each block look at near by blocks as well
- spatial trees just divide the space using a tree name explains it well
- why is lsh more efficient
- we are not partitioning and lowering the number of candidate pair

quiz questions

- min hash fails when a single element belongs to all sets in a collection
- true if any hash picks this item as a minimizer than all sets could

reproducibility

- reproducibility really important in big data
- there are all kinds of best practices
- reproducible folder standard
- keep contextual information in README
- sensitive data kept in secured repository

recommender systems

- idea predict which items a user will interact with
- method popularity model + dampening factor
- latent factor model model interactions as inner product of two vectors
- implicit vs explicit feedback

search , ranking evaluation

- can use min hash and LSH to identify similar documents
- PageRank orders documents by their probability of capture a random walk
- core computation is learning eigenvectors, and transition matrix
- can use power iteration

differential privacy

- being able to release data is critical for reproducibility but needs to be balanced
- k-anonymity is not enough
- de-anonymized attack highly accurate

- differential privacy works through an API
- add Laplacian noise to give plausible deniability
- sensitivity maximal difference in output given a single row different external aggregates high sensitivity each
- multiple queries reduce DP

GPUs

- computation parallelism dedicated hardware
- more restricted program control flow
- but less restricted in terms of data access
- limited sharing of information between computation threads
-