

title

wbg231

January 2023

1 Introduction

- the speed up from spark comes from delayed computing parallelism and caching
- it is not faster at every problem but in something like gradient descent it is way faster
- it is

column oriented storage

- it is all about speed (ie wall time)

history of column oriented storage

- idea goes back to the 80's
- had a resurgence in the 200s
- in the 80s it was not seen as that important because both cpu speed and storage speed were increasing
- starting in the mid 2000s storage speed stagnated so to make things faster speeds up were needed
- transferring from disk to memory is vary slow
- sequential memory reads are faster due to cache per fetching
- so we want to transfer fewer bites and use predictable and contiguous memory access patters to make reading data faster

row oriented storage csv

- if you had data stored in a row of text oriented way how would you go about getting the nth record, or just takes the kth column. you would have to go through each row until you find the kth value in that row since strings are of variable size
- rows and column are hard to predict
- basically requires a full serial scan

record oriented storage relational data

- relational data can be logically grouped by rows
- that is good if you want to process all records in a row at one time
- that is also nice for appending data
- it is human readable as well

queries row stores

- getting a col from a row oriented database is equivalent to a loop in the best case
- each row is loaded from storage
- attribute is inspected
- rows that pass are sent down stream
- an index can help locate rows but that still involves pulling entire rows when we only want one column
- loading data from the disk is slow

column oriented storage

- each column is stored on its own
- values in each col have a constant type
- disk access patterns become much more regular
- this improves locality
- enables compression and vectorized processing

Row-oriented	Column-oriented
id, Species, Era, Diet, Abundant	id: [1, 2, 3]
1, T. Rex, Cretaceous, Carnivore, True	Species: ["T.Rex", "Stegosaurus", "Ankylosaurus"]
2, Stegosaurus, Jurassic, Herbivore, True	Era: ["Cretaceous", "Jurassic", "Cretaceous"]
3, Ankylosaurus, Cretaceous, Herbivore, True	Diet: ["Carnivore", "Herbivore", "Herbivore"]
	Abundant: [True, True, False]

-
- i think the above picture shows why this works well as vectors.

speed is not everything

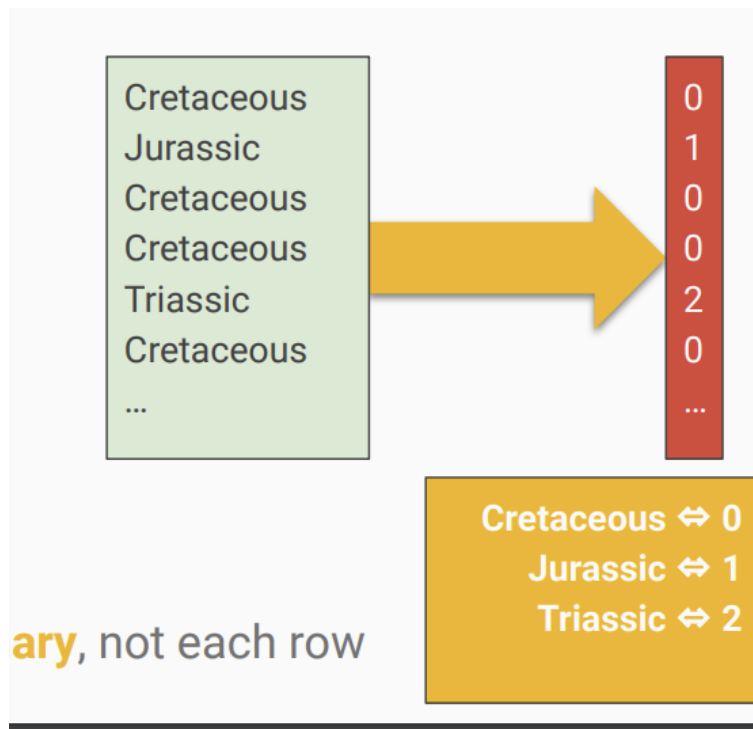
- storage space matters too
- mixed types are hard to compress
- once that data is arranged in a columnar fashion they all have the same type
- so they can be compressed saving space meaning that the data can also be sent and received much more easily

compression

- records have heterogeneous types
- a single column has one type
- that means there is low entropy in a column so can be easily compressed
- compressed data takes less space, is cheaper to load and sometimes we can compute directly on compressed data
- but what compression should we use

dictionary encoding

- this is where we encode each value in a column with a unique key
- this works well when we have a few distinct values
- replace string value by string identifiers this allows the column to have uniform data with and better cache locality
- so string matching can be done on the dictionary not each row (since they are stored in different dictionaries)



bit packing

- integers usually take 4-8 bytes to store (32 bits or 64 bits)
- bit packing squeezes small integers together

Values	0	1	0	2	1	1
8-bit (binary)	0000 0000	0000 0001	0000 0000	0000 0010	0000 0001	0000 0001
Compressed	0001 0010	0101				

- matching and comparing can be done on the compressed data
- this will only work well when there are a lot of small integers

run length encoding

- useful when we have long runs of constant values
- we convert a sequence of values to tuples of the type (value, #repetitions)
- sum, average, counts and other aggregations can all be done on compressed values

other compressions

- frame of reference encoding 1000, 1004, 1005, 1002 \rightarrow 1000|0, 4, 5, 2
- delta coding 1004, 1005, 1006, \rightarrow 1004| + 0, +1, +1
- many others
- compression can be combined
- the main trade off is space efficiency vs complexity of querying
- **slide** You work as a data scientist for Netflix and need to compress a movie to stream it efficiently. What is the most suitable compression scheme? Hint: A motion picture consists of many successive frames
- delta encoding, maybe run length

column oriented storage take away

- pros
 1. can be faster if we only want a subset of attributes
 2. higher storage efficiency and throughput
 3. collecting data of the same type allows for compression
- cons
 1. reconstructing full tuples from compressions can be slow
 2. writes and deletions can be slow
 3. handling non tabular data is tricky

2 when data is not tabular

dremel and parquet

dremel

- dremel is a low latency query system for read only **structured data**
- developed at google
- lots of cool ideas in the paper but let's talk about data format
- core ideas were quickly taken and reused in parquet

nested and structured data

- not everything fits nicely in relations
- variable lengths and depths are hard to deal with
- record oriented storage is more natural here
- how can we get all the benefits of column stores but for structured data

trees

- we use the hierarchical data structure tree

example web documents

- there are required and not required tags
- need a doc id
- don't need links need at least 1 name and 1 language execution

what specs would we like to see in a system that flattens hierarchical records

- we want lossless representations of the hierarchical records in columnar format
- it needs to be possible to recreate the hierarchical record from the columnar format
- the key challenge is being able to parse records unambiguously
- we need to be able to keep track of the record structure, ie if a value appears in a table row we need to understand whether it is the same piece of data or a unique record
- to make this efficient it needs to be able to handle sparse datasets
- we need to be able to represent missing fields efficiently like null values

implementing records flattening with dremel

- the key idea is keeping track of repetitions of fields within a record to parse
- the repetition level (which level repeated most recently) r
- the definition level d how many optimal diles are present
- the required fields same are the same level as the parents
- optimal fields the same r level as parents d level increments
- then there are repeated fields
- go back through the flattening example i don't have the life force for this right now
- dremel can easily rebuild partial views
- unused attributes can be ignored
- but decoding the data is sequential so dremel data is hard to parallelize

after flattening

- repetition and definition cols are highly compressible
- value fields are a new column of the same type
- cols broken into blocks and compressed independently

parquet

parquet

- developed at twitter in 2013
- default storage for spark
- based on dremel flattening but without the analysis engine or query machine

parquet format

- pages for a column are compressed independently
- small pages make it easier to read records but incur more overhead
- row groups should be large but fit into one hdfs blocks

nice things about parquet

- there is cross language support
- allows for partial decoding ie only look at a few cols
- it works well with spark and hdfs
- preserved rdd and data frame directly

using parquet in practice with spark

- column efficiency depend on row order it largely relies on how compressable the data is
- data frame partitions can be written out repeatedly
- most frameworks we use in coding are already column oriented