

# Big data in class final

wbg231

January 2023

## Introduction

- ok they are getting started
- final is friday at 6pm
- need pencil and id
- can bring 1 page of hand written notes
- don't bring any devices
- exam consists of 60 multiple choice questions which will be evenly distributed
- there are only conceptual questions, not code. this is stuff that could come up in an interview high level conceptual questions
- just hand write notes, lets just keep stuff simple aim to write less.
- going to be on a scan tron
- focus on things that are to some extent at a high level
- around the same level of detail as the quiz

## survey results

- there were not many survey results
- tried to make review slides follow that distributions
- most were focused near the 3rd quarter of the class.

## **relational databases**

- not much to say
- they help standardize how data is
- data validation is also helpful
- speed up access
- sql is good
- ACID rules
- atomicity if have a sequence of operations either all complete together or all fail together
- related to relational databases but not specific to it, it is a thing for any concurrent databases
- durability is where we pass to the operating system
- consistency just means that everything will stay in line with the schema always in valid started
- independence order of execution is irrelevant this matters less because we have read only data in most stuff

## **map reduce**

- map: process 1 record at a time generate key value pairs
- reduce: process 1 key at a time output single object (these are parallel)
- combiner: kind of like a reducer in the mapper
- constraining form of computation to make things faster (really really tight constraints )
- combiner tries to minimize key skew, by doing a partial reduction
- combiners only really work if certain conditions are met

## **complexity analysis**

- analyzing complexity
- just count the number of loops my man sheesh
- input of some length

- how many steps does that algorithm take to complete that results
- we only care about the dominating term and the lower bound  $O(n)$
- can do the same thing for space
- in the context of map reduce have to loop at the number of keys, number of reducers, number of mappers
- questions about combiner,

### **map vs reduce**

- reducer always sees key and list of values for that key, mapper just sees input keys
- reducer sees all the values for a key, mapper just sees one record at a time
- the mapper should be independent of one another

### **map reduce questions**

- reducers do not reduce the number of keys, they reduce the number of values per key the number of jobs that goes to the reducer will be the same
- a map function in map reduce does not have to produce any intermediate value for every line, like if we were filtering for a word

### **HDFS**

- name nodes vs data nodes
- data nodes store the data
- name nodes store the name
- in HDFS large file broken into blocks, the data nodes store the blocks but not the meta data saying where the data is from
- the name node maps file names to the data nodes that contain the block, the name node is your lookup table the data node is your storage
- in HDFS we get around concurrent access stuff, by just keeping the files read only
- we do not have to worry about t

- replication factor, default replication rate is 3. 1 copy on 1 machine, 1 copy on a different machine on the same rack, 2nd copy on a separate rack if that rack goes down
- raising the replication factor means that it is going to be easier to find the data that you need at the cost of storage
- does it cost computation efficiency to increase replication rate, no it is increasing data? it will increase storage cost obviously
- data is append only, does that mean we can add info to the same file?, the idea is that you are allowed to append ie just change the last block in the file once the block is full you create a new block.
- this is why when you run the same map reduce program without cleaning up output files it crashes?
- the mapper writes back to temp files on hdfs that are cleaned up by map reduce,

## questions

- if the name node fails you are in trouble.
- it happens it is not permanent it just means there is server down time
- this goes back to the CAP theorem, we are giving up A when the name node fails, you need to wait for the server to not be crashing

## spark

- the bottom line is that map reduce is from mid 2k it was good for what it was meant to do but not all tasks
- but we wanted to do other tasks, like iterative tasks,
- in map reduce all the reducers sit around idly until the mappers finish that is not good if we are doing iterative stuff
- spark uses RDD where we separate actions from transformations. we can do a lot of transformations this allows for efficient computation
- **wide vs narrow dependency** a narrow dependency is easy to parallelize. the solution is almost always make sure that if you are joining RDD's they have the same partition structure (and thus are a narrow dependency) the partitions can map directly, the output depends on at most one input partition so can propagate the same partition structure

- a wide partition is something that depends on all our data, joins are wide unless two tables have same partition structures
- want dependency narrow or it is hard to parallelize
- spark is a better implementation for SQL as well as machine learning
- SPARK gets the same parallelism as map reduce with less restrictions
- spark is a direct response to stone breaker
- relation between spark and map reduce, spark is separate from map reduce it is a completely separate thing
- ...

### **spark questions**

- compare and contrast rdd and data frame
- spark data frames are not read only that is the whole point that is why we have iterative stuff going on.
- each step in an rdd lineage graph must be completed before the next false
- spark uses pipelines to enclose multiple stages of map reduce processing. false no it is a separate thing that was made to address issues from map reduce
- spark can use hadoop but hadoop and map reduce are not the same thing
- are the actions going to be sequential,
- lineage graph is just about transformations, the actions are outside of the rdd
- rdd's move backwards through the computation graph until it finds a data source

### **column oriented storage**

#### **review**

- the default for spark is parquet
- this makes things faster sometimes just need 1 column but not all rows.
- regularity helps with that this leads to speed ups
- column oriented storage because columns have same data types all us to search through them faster

- can use dremel that always turns a hierarchical object into a tabular representation that we can view as a column oriented storage
- the bigger take away is that dremel was cool but no one uses it
- parquet is actually used.
- csv files are a lot slower than parquet
- this is why the project data was all parquet not csv
- parquet is not human readable though

## questions

- explain parquet and dremel, just understand the relationship between the two .
- the dremel system was designed to process all subsets of records in a data set. in it is for subsets of attributes for all records. the point is it is not record based it is feature based
- when written to hdfs parquet files localize different columns in different hdfs blocks that is false, parquet files divide blocks by rows not column
- what is a block is a subset of records a contiguous chunk of a dataframe but spanning all the columns, that is why the whole thing is called parquet

## dask

### review

- there were a lot of dask thoughts
- when would you want to use dask?
- it is kinda dependent, spark is good if you are ok doing everything in spark if you need to do things where spark interacts with something in python that is not written in spark, then it gets a bit messy to put things in and out of spark. so dask is nice for that kinda thing
- most of the time you will be dealing with sql
- after that spark number 2, spark is good for regular data
- for ragged data or messy data dask can be useful, or if you need parallelism that is not rows in a data frame but chunks in an array, get more flexibility in your data types

- scalal makes constrains makes stuff faster
- dask accepts what ever comes out of your pyhton code, but hard to figure out waht you are doing
- you have to work harder to make dask work well
- but there are some jobs where dask is ideal
- daks is for data scicne use cases not software engeneering stuff

## questions

- cross compare spark and dask. dask is spark written by scipy poeple, the computational principles are similar. in terms of how to optimize stuff it comes back down to partiion sturcutre need to minimize comuncaiton early on as well
- in some programs just have to loop, the thing is you need to identify when the task is worht it.
- no strong words of advice
- in sql indecies are how you optimze those should be designed for the queries you think you want to run
- those are just additional datastrucutres that can make things faster but do not effect ur data

## aproximate nearest neighbors

- this is all similarty search
- one of the key aplications is sim saerch
- that is a fundemtnal use case
- as you increawse block size in lsh your liklyhood of colision goes down but have to do more work, so there a trade off with block size and number of blocks
- similarty search lsh lets you do this in a really fast way
- boost the high similarty item reduce the low similarty item
- brute force similarty is overwhelming for large documnets that is the high lvl thing.
- we did not really talk about distance metrics

## questions

- lsh is built on the idea of randomness in mini has that is from the hash function but that can be wasteful spatial tree is instead of doing independent partitions we do recursive splitting, we are cutting our data point in half as many times as we want. this is effectively many locally dependent hash functions
- cosine similarity lsh if you want to approximate cosine similarity can check the lsh by randomly projecting them
- multi probe not requiring that strict a match for all values
- multiprobe at every state fuzz the results a bit
- at the high level cosine similarity is very general
- min has failures when a single element belongs to every set in a collection true
- min has signatures are generated by applying false the min hash functions has not directly formed in min hash

## reproducibility

- this is really important in big data
- it is hard to reproduce results in big data and they propagate over large dataset
- reproducibility can you reproduce this.
- most of the results from the 90's can not be reproduced
- there are all kinds of best practices that help with reproducibility
- want to have three copies of all your data files, at least 1 backup off site?
- want to have organized project structure, input data, process code, results results, metadata
- want to have at least one readme file on how to use things
- in jupyter notebooks some times cells are not even run in order
- there are safer ways to store sensitive data
- version control use git
- that is super important in big data
- get the same result so can retrace your steps



## recomender systems

- idea is to predict which items a user will interact with from a large catalog
- what is the dampening factor? it is how we considered items with more interactions to be a better estimate of interactions and make sure they weigh more
- the next step is the latent factors model
- decompose the model into separate vectors
- implicit vs explicit feedback, explicit like ratings implicit like if you watch a youtube video or watch ever
- explicit feedback is usually a stronger signal that is a clear indicator but is often of a lower volume
- implicit feedback tends to higher volume with lower signal to noise feedback,

## searching ranking evaluation

- google used random surfer
- page rank etc
- ranked lists evaluations for search recommendation systems etc
- the thing that we say is that this cooks down to a linear algebra problem
- rank list evaluation is important
- the simplest model is that we make a predicted list of items for each user, we want all the positives to come before all the negatives. the details between how you weight positives and negatives and where they are in the list distinguished metrics, but the main point is evaluating a rank list as if it is a prediction

## question

- in page rank, the model is random if that browser lands on a page that has no outgoing links or links to itself, then the model can not get out of there
- so if you ever have a change of landing there, it will think it is important, so they introduce this teleportation parameter that lets your surfer teleport to somewhere random

## differential privacy

- ways to keep data more private
- in one sentence it is shown that just anonymity is not sufficient
- guess who you can reduce each space in a lot of dimensions
- once the thing is de-anonymized it is out there :/ if that is sensitive we have problems
- differential privacy keeps raw data private for ever and you interact with the dataset through an api
- add some specific laplacian noise that can not easily be removed it is absolute value of exponential so hard to remove
- it is hard to know if you are in the data set if it is given
- for a larger dataset there is a lot of plausible deniability
- this is limited to multiple queries
- so privacy laws agree
- the noise that we inject is at the level of the computation we ask in result for it is on a per query basis not, in the dataset it's self
- the raw data is unaffected which is nice :)
- but we inject noise on the level of the query so it allows for plausible deniability, the user of the api may know what type of noise was put in the dataset

## gpu

- computational parallelism from dedicated hardware
- similar to traditional cpu but with more restrictive program flow
- same components just in different proportions
- gpu many small processors that do a few computations in small steps
- think of each core on the gpu as touching a different part of memory that is a lot of the way there
- cuda framework is thread ie individual computations grouped into blocks which are grouped into grids
- and that is how we divide the work up in a way that we can divide over our entire dataset we just run on one grid at a time until we are done

- can return arbitrary data types
- a cuda kernel does not prevent you from looking at different data
- that is nice because you can do convolutions easily
- this allows limited communication between threads but doing this effectively requires some involved programming
- they are starting to creep more into like normal data science applications
-