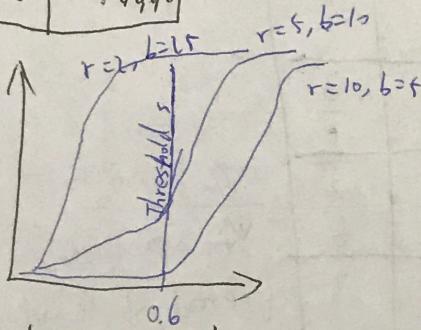


S	$1 - (1-S)^b$	at least one bond is identical
0.2	0.036	
0.3	0.071	
0.4	0.186	
0.5	0.470	
0.6	0.812	
0.7	0.975	
0.8	0.996	



\Rightarrow Tune M, b, r, t , get
 (almost) all pairs with similar signatures.
 Delimit more pairs that do not have similarity.

element X_{n+1} arrives what is the probability that an already-sampled element X_i stays in the sample

$$(1 - \frac{S}{n+1}) + (\frac{S}{n+1}) \cdot (\frac{n-1}{n}) = \frac{n}{n+1}$$

X_{n+1} not sampled

X_{n+1} sampled

X_i not evicted

Hence, in the sample with probability

$$\frac{S}{n+1} \text{ chance}$$

(1) \rightarrow keep track of sampled items

(2) $O(N)$

(3) $O(n^2)$, n is the number of items

Data Stream
Sampling a fixed proportion
固定比例采样

$$\frac{S}{n} \Rightarrow S \text{ Sample size}$$

存储缓冲区

Reservoir Sampling

Elements X_1, \dots, X_n

1. Store all first elements $X_1 \dots X_S$
2. Suppose element X_n arrives with probability $(1-S/n)$, ignore this element with probability S/n ;
 → Discard a random element from the reservoir
 → Insert element X_n into the reservoir

X_{n+1} is sampled with probability $S/(n+1)$
 If it is sampled $\rightarrow S/n$
 stay in the sample with probability $[n/(n+1)]$ 已不在缓冲区

Morris Algorithm

$x \leq 0$

for each of the n entries

$$x \leftarrow x + 1 \text{ with } (1/2)^r$$

$$\text{Return } n' = 2^r + 1$$

Counter x needs only $\log(n)$ bits

Let $X(n)$ denote count after

arrival n . Transition $x \rightarrow x+1$

with probability $(1/2)^r$

$0 \rightarrow 1 \rightarrow 2 \rightarrow \dots \rightarrow 2^r \rightarrow 2^{r+1} \rightarrow \dots \rightarrow X$

Universality

initialize $x=0$,
 increment w.p. $p=2^{-r}$,
 estimate $n'=2^r + 1$

$n=1$
 before: $x=1$, $p_0=1$
 Prob. $1: X \rightarrow 1$
 estimator $n'=2^1 + 1 = 3$

$n=2$
 before: $x=1$, $p_1=1/2$
 prob. $1/2: X$ stays at 1
 $n'=2^1 + 1 = 3$
 Prob. $1/2: X \rightarrow 2$; $n'=2^2 - 1 = 3$

Hierarchical Clustering (Top-Down) Dividing k -Means

For $I=1$ to $I \leftarrow k$ do

Pick a leaf cluster C to split

For $J=1$ to $ITER$ do

use k -means to split C into two sub-clusters C_1 and C_2

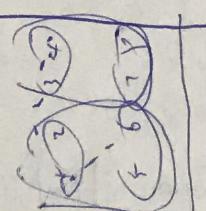
choose the best of the above splits and make it permanent

Bottom Up Clustering

key operation: repeatedly combine two nearest clusters.

3 important operations:

① How to represent a cluster of more than one point?



② How do you determine the 'nearness' of clusters?

③ When to stop combining clusters

clear

single

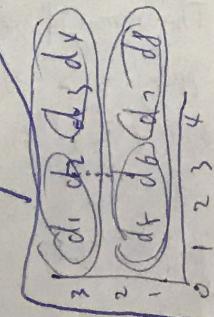
$$\text{Sim}(C_i, C_j) = \max_{x \in C_i, y \in C_j} \text{sim}(x, y)$$

$$\text{Sim}(C_i \cup C_j, C_k) = \max(\text{sim}(C_i, C_k), \text{sim}(C_j, C_k))$$

Complete Link

$$\text{Sim}(C_i, C_j) = \min_{x \in C_i, y \in C_j} \text{sim}(x, y)$$

$$\text{Sim}(C_i \cup C_j, C_k) = \min(\text{sim}(C_i, C_k), \text{sim}(C_j, C_k))$$



Mid-term

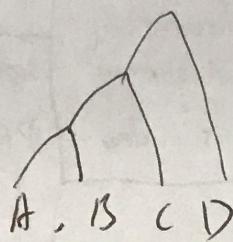
(Rule with 100% confidence are not interesting
people more to have nothing
⑦)

Movie IP

	A	B	C	D	E
1	ABE	1	1	1	1
2	AP	2	1	0	1
3	BC	3	1	1	0
4	ABD	4	1	1	1
5	AC	5	1	1	0
6	BC	6	1	1	0
7	AC	7	1	1	0
8	ABCDE	8	1	1	1
9	ABC	9	1	1	1

	A	B	C	D	E
1	1	1	1	1	1
2	2	1	0	1	1
3	3	1	1	0	0
4	4	1	1	1	1
5	5	1	1	1	0
6	6	1	1	0	0
7	7	1	1	0	0
8	8	1	1	1	1
9	9	1	1	1	1

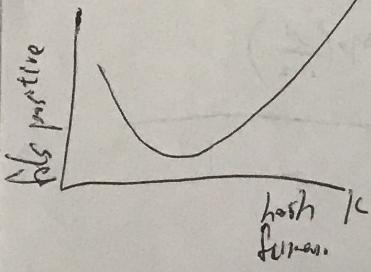
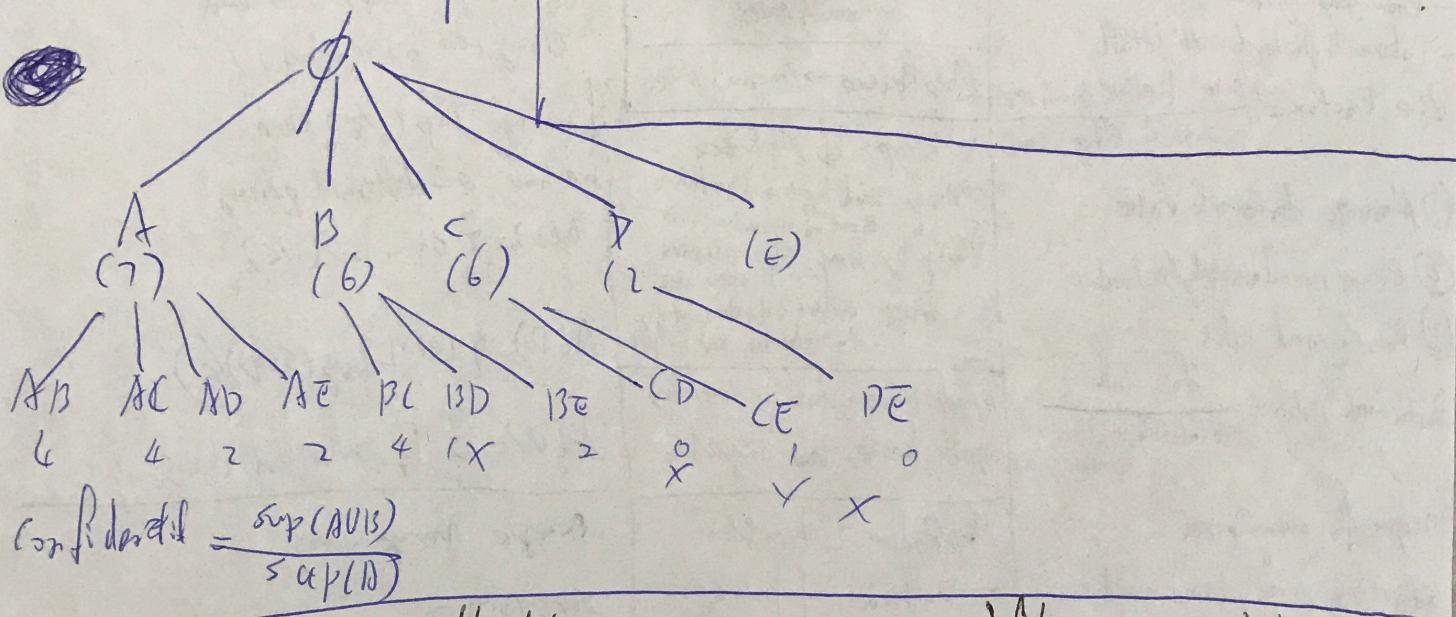
ABCD



	A	B	C	D
A	0	1	3	7
B	1	0	2	6
C	2	1	0	4
D	3	2	4	0

$$d(AB, CD) = \min(AC, BC)$$

	P	A	C	D
P	0	1	0	0
A	1	0	1	0
C	0	1	0	0



handle deletion

use counters instead of 0/1s

when adding an element increase the covers

when deleting an element, decrease the covers

covers must be large enough to avoid overflow (kb-r)

$$\frac{1 - (1 - 1/n)^{n(m/n)}}{1 - e^{-m/n}} < 1 - e^{-m/n}$$

some target X not hit by a dom

prob of large one does hit X

false positive probability fp
 $e^{-km/n}$

Bloom filter operators

Insert(lc)

for i from 1 to r

$$B[h-i(lc)] \leftarrow 1$$

CSMember(lc)

for i from 1 to r

$$if B[h-i(lc)] == 1$$

return false

return true

why visualization

- ① Record
- ② Analyze data to support reasoning
- ③ communicate

Complex ideas communicated with clarity, precision, and efficiency

substance, statistics, Design
Tufte's Principles

Class Design Principles:

- no cute, on focus data
- have a good explanatory, clear title
- Data integrity; maintain scale, right axes
- know your goal
- Balance high-level with detail

$$\text{Lie Factor} = \frac{\text{size of effect shown in graph}}{\text{size of effect in data}}$$

① Maximize data-ink ratio

② Erase non-data-ink/redundant

③ Revise and edit

$$\text{Data-ink ratio} = \frac{\text{data-ink}}{\text{Total ink used to print graph}}$$

argue for every pixel

Starting point: erase non-data ink
erase redundant ink

encoding data with shapes, color, and size. Which cues to choose depends on your data and your goals

logarithmic $\rightarrow \log \rightarrow$ percent change

Correlation / distribution

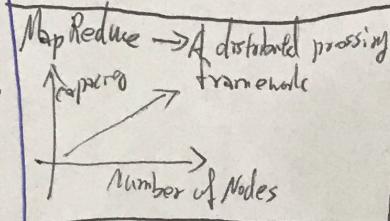
Clustering / outliers

Good story

introduce interesting characters
challenge is believable
There are hurdles to overcome
the outcome or prognosis is clear

Map Reduce

The Hadoop Distributed File System (HDFS)



MapReduce \rightarrow Pragmatic Model

- 3 steps of MapReduce
- ① Map: Read input and produce a set of key-value pairs
- ② Group by key: collect all pairs with same key
- ③ Reduce: collect all values belonging to the key and output

spark process data in-memory

spark better than mapReduce

MapReduce 2 limitations

- Difficulty of programming directly in MR
- performance bottleneck, or both not fully.

Spark \rightarrow RDD

Resilient Distributed Dataset (RDD)

- across the cluster, Read-only
- caching datasets in memory

RDD \rightarrow part of Go's spark

Amazon Pyramid

$-t \rightarrow \text{IPF} \rightarrow \text{Amazon} \rightarrow \text{replica}$

Sensitivity

$$\text{PDF: } f(x|\mu, b) = \frac{1}{\sqrt{2\pi b^2}} \exp\left(-\frac{(x-\mu)^2}{2b^2}\right)$$

mean μ variance $2b^2$

$$|q(D) - q(D')| \leq s(q)$$

Example: Count query

Number of people having disease

sensitivity = 1

"differing in one row" \rightarrow add change count by max of 1

if we use lap(1/ε) then

we have ε-differential privacy

$$A(D) = \sum_{i=1}^n d_i + \text{lap}(1/\epsilon)$$

$$A(D) = q(D) + \text{lap}(s(q)/\epsilon)$$

$$f(v) = \sum_{i=1}^n d_i \text{ where each } d_i \in S_{0,1}$$

example: Average query

sensitivity =

$$M/n$$

$$A(D) = \frac{1}{n} \sum_{i=1}^n d_i + \text{lap}\left(\frac{m}{n\epsilon}\right)$$

Statistical Data Privacy

目标：个人信息被保护

使用目标：统计数据有用

① Anonymization → linkage attacks

② Query auditing → 确保查询结果准确

③ Summary statistics → differentially private

可以 offline/online 来检测
 回答 request → 看是否回答 request
 带有噪声的 response

DIFFERENTIAL PRIVACY

不论是否加入了 db，统计结果无法被推断
 (perturbed with differential privacy)

$$\log \left(\frac{\Pr[A(D_1) = 0]}{\Pr[A(D_2) = 0]} \right) \leq \epsilon \quad (\epsilon > 0)$$

不论什么 record 加入，都不变

$$\Pr[A(D_1) = 0] \leq e^{\epsilon} \Pr[A(D_2) = 0]$$

控制 D_1, D_2 的分离程度。

ε 越小，越 privacy

$$\frac{\Pr[A(D_1) \in Y]}{\Pr[A(D_2) \in Y]} \leq \exp(\epsilon)$$

$\exp(-\epsilon)$ 对于 ε 很小的情况
 $1 - \epsilon \leq \sim \leq 1 + \epsilon$

$$\text{Laplace: } f(x | \mu, b) = \frac{1}{2b} \exp\left(-\frac{|x-\mu|}{b}\right)$$

Theorem: if sensitivity of the query is S , then the algorithm $A(D) = q(D) + \text{Lap}(S(q)/\epsilon)$ 确保 ϵ -differentially private

如果算 COUNT → $b - t$ 值，则

$$A(D) = \frac{n}{f(D)} \sum_{i=1}^n d_i + \text{Lap}(1/\epsilon)$$

$$\text{Sensitivity} = \frac{1}{n}, A(D) = \frac{1}{n} \sum_{i=1}^n d_i + \text{Lap}\left(\frac{1}{n\epsilon}\right)$$

M_1, M_2 都是 D 的某 $\frac{1}{k}$
 满足 ϵ -differential privacy

那么，
 结果： $\hat{E} = E_1 + E_2 + \dots + E_k$

Sequential Composition

M_1, M_2 是 D, \dots, D_k 的输出
 为 E_1, \dots, E_k , 且 ϵ -differential
 $\max\{E_1, \dots, E_k\}$

$M_1 \rightarrow D$ 为 ϵ differentially private

$M_2(M_1(D))$ 也满足 ϵ

K-Mean with Differential Privacy
 $\epsilon/T \rightarrow \epsilon$ after T iterations

1. 分配点
2. 随机扰动簇中心
3. 计算 noisy sum

$$h(x) = x \bmod S$$

$$g(x) = (2x+1) \bmod 5$$

$$\text{Sim}(c_1, c_2) = |C_1 \cap C_2| / |C_1 \cup C_2|$$

Jaccard distance

$$= 1 - \text{Sim}$$

无混淆 bucket: ① 相同样本
 ② 不同文件
 ③ 不同文件

documents as sets of shingles

why shingle? set of tokens
 ① 保证不重复
 ② 保证大小写
 ③ 保证不分割

为什么用 shingle? 10% error rate

$$N(N-1)/2 \approx S^2 / 10^6$$

把 docs 放进 bucket, filter by
 docs 会放进 different bucket

random permutation π

$$h_\pi(c) = \min_m \pi(c)$$

Similarity

Single-pass method

for each row r

for each hash function h_i
 compute $h_i(r)$

for each column c

if c has 1 in row r
 for each hash function h_i
 if $h_i(r) < M(t, c)$
 then $M(t, c) < h_i(r)$

Row	C ₁	C ₂	h(x)	g(x)
1	△1	0	1	3
2	0	△1	2	0
3	△1	1	3	2
4	△1	0	4	4
5	0	△1	0	1

Same bucket means
 "identical" in the band

① Find Pair of $\geq s = 0.8$
 $b = 20, r = 5$

② Probability C_1, C_2 identical in
 one particular band: $(0.8)^5$
 —→ Band 完全相等 = 0.328
 ③ all robust not similar
 $(1 - 0.328)^{20} = 0.22035$
 完全随机的几率 ↑

1/30000 2 false negative
 Find 99.961% pair true
 相似

but $\text{Sim}(c_1, c_2) = 0.3$
 $0.3 < 0.8$

+ 5% 例从桶中取后需要
 找到 common bucket

$(0.3)^5 = 0.00243$
 $(1 - 0.00243)^{20} = 0.99757$

4.7% candidate pairs
 中等概率
 ↓ False positives

if rebound only if bands/s
 rows, false positive ↓,
 false negative ↑

all rows in band equal = 1
 ② Some row in band unequal =
 $1 - t^r$

③ no band identical =
 $(1 - t^r)^b$

④ 1 band identical
 $1 - (1 - t^r)^b$

每份 hash 表 ≠ hash table
 with k bucket. 1c 表大表

1/b 和 1/c 表示相等 pair