

drawback

- does not understand similar meaning words.
- no semantic similarities between words.

Pearson correlation

	I ₁	I ₂	I ₃	I ₄	I ₅
Alice	5	3	4	4	?
User1	3	1	2	3	3
User2	4	3	4	3	5
User3	3	3	1	5	4
User4	1	5	5	2	1

$$\text{Sim}(a, b) = \frac{\sum_{p \in P} (g_a, p - \bar{g}_a)(g_b, p - \bar{g}_b)}{\sqrt{\sum_{p \in P} (g_a, p - \bar{g}_a)^2} \sqrt{\sum_{p \in P} (g_b, p - \bar{g}_b)^2}}$$

$$\text{Sim}(\text{Alice}, \text{User1}) = \frac{(5-4) \times (3-2.4) + (3-4) \times (1-2.4) + (4-4) \times (2-2.4) + (4-4) \times (3-2.4)}{\sqrt{(5-4)^2 + (3-4)^2 + (4-4)^2 + (4-4)^2} \sqrt{(3-2.4)^2 + (1-2.4)^2 + (2-2.4)^2 + (3-2.4)^2}}$$

$$= \frac{2}{\sqrt{2} \sqrt{2}}$$

$$= 0.839$$

$$\text{Sim}(\text{Alice}, \text{User2}) = \frac{(5-4) \times (4-3.8) + (3-4) \times (3-3.8) + (4-4) \times (4-3.8) + (4-4) \times (3-3.8)}{\sqrt{(5-4)^2 + (3-4)^2 + (4-4)^2 + (4-4)^2} \sqrt{(4-3.8)^2 + (3-3.8)^2 + (4-3.8)^2 + (3-3.8)^2}}$$

$$= \frac{1}{\sqrt{2} \sqrt{1.36}}$$

$$= \frac{1}{1.649}$$

$$= 0.606$$

$$\text{Sim}(\text{Alice}, \text{user } 3) = \frac{(5-4)(3-3.2) + (3-4)(3-3.2) + (4-4)(1-3.2) + (4-4)(5-3.2)}{\sqrt{(5-4)^2 + (3-4)^2 + (4-4)^2 + (4-4)^2} \sqrt{(3-3.2)^2 + (3-3.2)^2 + (1-3.2)^2 + (5-3.2)^2}}$$

$$= \frac{-0.2 + 0.2 + 0 + 0}{\sqrt{2} \sqrt{13.56}}$$

$$= 0$$

$$\text{Sim}(\text{Alice}, \text{user } 4) = \frac{(5-4)(1-2.8) + (3-4)(5-2.8) + (4-4)(5-2.8) + (4-4)(2-2.8)}{\sqrt{(5-4)^2 + (3-4)^2 + (4-4)^2 + (4-4)^2} \sqrt{(1-2.8)^2 + (5-2.8)^2 + (5-2.8)^2 + (2-2.8)^2}}$$

$$= \frac{-1.8 - 2.2 + 0 + 0}{\sqrt{2} \sqrt{13.56}}$$

$$= \frac{-3.96}{5.208}$$

$$= 0.760$$

most similar users are user 1, user 4

~~unknown rating between Alice~~

$$r_{\text{Alice}, I_5} = \frac{4 + 0.839(3-2.8) + 0.760(1-2.8)}{0.839 + 0.760}$$

$$= 4 + -2.181 +$$

$$= 4 + \frac{0.503 - 1.368}{1.599}$$

$$= 4 - 0.541$$

$$= 3.459$$

Association rule mining

$X \rightarrow Y$
 Antecedent Consequent

$$\text{confidence} = \frac{\text{support}(X, Y)}{\text{support}(X)}$$

$$\text{lift} = \frac{\text{support}(X, Y)}{\text{support}(X) \times \text{support}(Y)}$$

$$\text{leverage} = \frac{\text{support}(X, Y) - \text{support}(X) \cdot \text{support}(Y)}{\text{support}(X) \cdot \text{support}(Y)}$$

assuming 1000 transactions with milk, egg appearing in 300 of them, milk 500, eggs 400, bread 400, milk, bread 400
 Find lift (milk \rightarrow egg) = ?
 lift (milk \rightarrow bread) = ?

If lift value is 1 X and Y are independent of each other

- lift ($X \rightarrow Y$) ≥ 1 indicates that there is a positive relationship
- Larger value of lift indicates the stronger association relationship between X and Y

Leverage

If leverage is 0, X and Y statistically independent of each other.

> 0 , X and Y have a relationship

a larger value of leverage indicates ^{relationship} stronger association rules between X and Y .

Inference

lift = 0 $\Rightarrow X, Y$ are independent

$> 1 \Rightarrow$ +ve correlation

$< 1 \Rightarrow$ -ve correlation

$$\text{lift}(\text{milk} \rightarrow \text{egg}) = \frac{\text{support}(\text{milk}, \text{egg})}{\text{support}(\text{milk}) \times \text{support}(\text{egg})}$$

$$= \frac{\cancel{300}}{\cancel{400} \times 500} = \frac{0.3}{0.5 \times 0.4}$$

$$= 0.0015 = 1.5$$

$$\text{lift}(\text{milk} \rightarrow \text{bread}) = \frac{\text{support}(\text{milk}, \text{bread})}{\text{support}(\text{milk}) \text{support}(\text{bread})}$$

$$= \frac{\cancel{100}}{\cancel{500} \times 400} = \frac{0.4}{0.5 \times 0.4}$$

$$= 0.002 = 2$$

$$\text{support}(x) = \frac{\text{no. of times } x \text{ occurs}}{\text{no. of transactions}}$$

$$\text{leverage}(\text{milk} \rightarrow \text{egg}) = \frac{\text{support}(\text{milk}, \text{egg})}{\text{support}(\text{milk}) \text{support}(\text{egg})}$$

$$= 0.3 - 0.5 \times 0.4$$

$$= 0.1$$

$$\text{leverage}(\text{milk} \rightarrow \text{bread}) = \frac{\text{support}(\text{milk}, \text{bread})}{\text{support}(\text{milk}) \text{support}(\text{bread})}$$

$$= 0.4 - 0.5 \times 0.4$$

$$= 0.2$$

\Rightarrow milk \rightarrow bread is strong association.

Streaming Data

26/5/22

e.g.: Google Search Queries in given period of time

Query

- data arrives continuously and rapidly
- user don't have control
- not uniformly distributed.
- cannot be stored in server
- * sampling operations
- * analytics is performed and summary is stored.

outputs

- approximation / estimation
- probability of max value.

Sliding window

Applications

- counting the distinct elements present in the stream (Flajolet Martin algorithm)
- filtering / selecting data (Bloom's filter)
- count the number of one's in the window (GIMM Algorithm)

Flajolet Martin Algorithm (FM)

- Selecting a Hash function h such that each element in the set is mapped to a string of atleast $\log_2 n$ bits.
- For each element $r(x) = \text{record}(x) = \text{length of trailing zeros in } h(x)$.

$$R = \max(r(x))$$

$$\begin{array}{l} \text{distinct} \\ \text{no. of elements} \end{array} = 2^R$$

Find no. of distinct elements

1, 4, 2, 1, 2, 4, 4, 4, 1, 2, 4, 1, 7

Hash function $h(x) = 3x + 1 \bmod 5$

$$\begin{aligned} 1 &\Rightarrow h(1) = 3(1) + 1 \bmod 5 \\ &= 4 \bmod 5 \\ &= 4 \end{aligned}$$

$$\begin{aligned} 4 &\Rightarrow h(4) = 3(4) + 1 \bmod 5 \\ &= 13 \bmod 5 \\ &= 3 \end{aligned}$$

$$\begin{aligned} 2 &\Rightarrow h(2) = 3(2) + 1 \bmod 5 \\ &= 7 \bmod 5 \\ &= 2 \end{aligned}$$

$$\begin{aligned} 1 &\Rightarrow h(1) = 4 \\ 2 &\Rightarrow h(2) = 2 \end{aligned}$$

$$4 \Rightarrow 3$$

$$4 \Rightarrow 3$$

$$4 \Rightarrow 3$$

$$1 \Rightarrow 4$$

$$2 \Rightarrow 2$$

$$4 \Rightarrow 3$$

$$1 \Rightarrow 4$$

$$7 \Rightarrow h(7) = 3(7) + 1 \bmod 5 = 2$$

∴ 4, 3, 2, 4, 2, 3, 3, 3, 4, 2, 3, 4, 2 ~~2~~

Convert to binary:

100, 011, 010, 100, 010, 011, 011, 011, 100, 010, 011,
100, 010

Find the trailing no. of zeros

2, 0, 1, 2, 1, 0, 0, 0, 2, 1, 0, 2, 1

record maximum value = 2

distinct number of elements = $2^R = 2^2 = 4$

DGIM algorithm

- provides 50% approximate output.
- Datar Grinias Indyak Motwani
- count the number of 1's in the window.
- rules that must be followed when representing a stream by bucket
 - right end of a bucket is always a position with 1
- every position with a one is in some bucket
- no position is in more than one bucket
- There are one or two buckets of any given size upto some maximum size.
- all sizes must be a power of 2.
- buckets cannot decrease in size as we move to the left.
- estimate the no. of 1's in the window with an error of no more than 50%. Purchase of items is represented as 1.

1 0 1 0 1 0 1 1 1 0 0 0 1 1 1 1

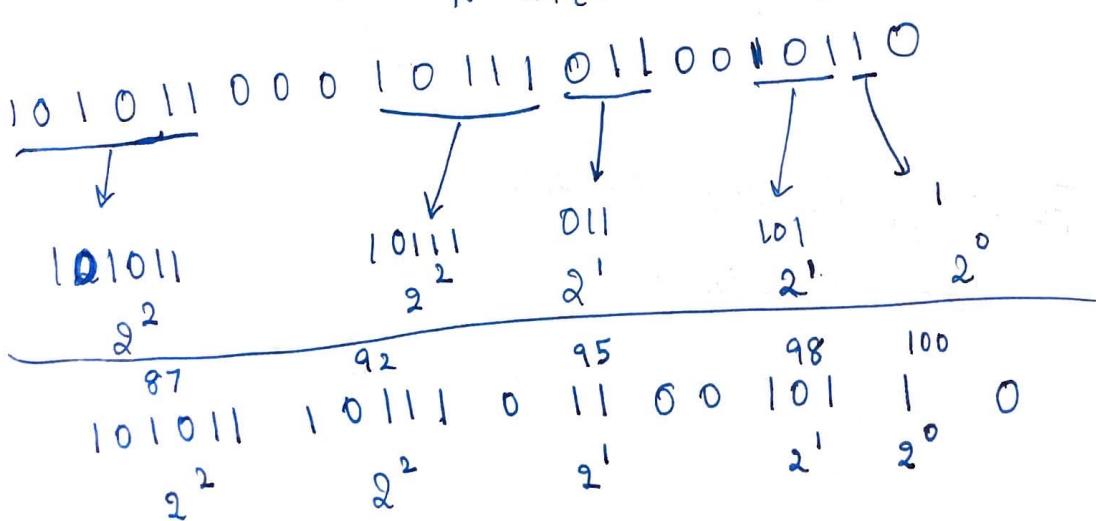
where 1 represents purchase of an item.

→ we divide the window into buckets the timestamp of its rightmost mention the most recent end.

- No. of one's in the buckets, the number must be a power of 2.

31/5/22

$$N = 24 \text{ [Window Size]}$$



new bits

← 101 102 103 104
0 1 1 1

At timestamp 101,

1010 10111 0 11 00 101 10
102 103 104

new bit zero enters

101011 000 10111 0 11 00 101 1 0 0 ← 102 103 104
 2^2 2^2 2^1 2^1 2^0 2^1 2^0 1 1 1

new bit 1 enters

101011 000 10111 0 11 00 101 1 00 1 1
 2^2 2^2 2^1 2^1 2^0 2^0 2^0 104
1

101011 000 10111 0 11 00 101

100	101	102	103
1	0	0	1
2 ⁰	2 ⁰	2 ⁰	2 ⁰

 2^2 2^2 2^1 2^1 2^0 2^0

101011 000 10111 0

11	00	101
2 ¹	2 ¹	2 ¹

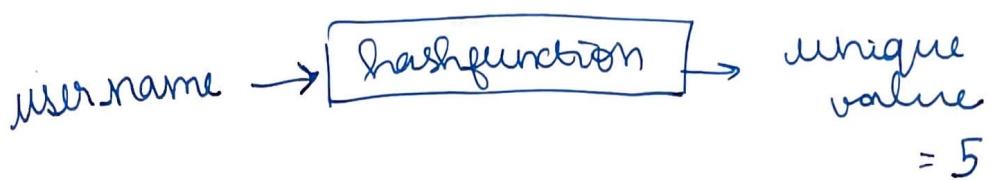
 1001 103
 2^2 2^2 2^1 2^1 2^0 1

Filtering / selecting stream

Bloom filtering

- always False negatives \otimes are not given but ~~long~~ there are probability for false positive

- works based on hashing technique -

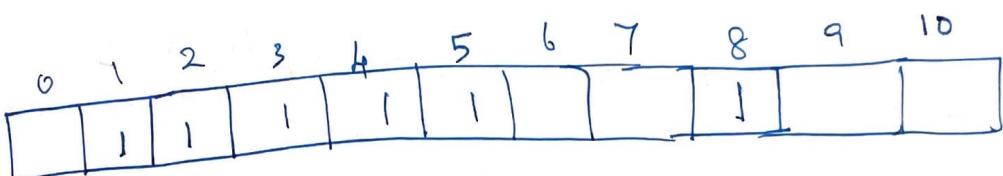


- to avoid collision, no. of hash function is ~~is~~ increased.

Insertion

1. ceg \Rightarrow $h_1(\text{ceg}) = 3$
 $h_2(\text{ceg}) = 4$
 $h_3(\text{ceg}) = 5$

2. mit \Rightarrow $h_1(\text{mit}) = 8$
 $h_2(\text{mit}) = 2$
 $h_3(\text{mit}) = 1$



Searching an element

ceg $h_1(\text{ceg}) = 3 \Rightarrow 1$ ————— AND
 $h_2(\text{ceg}) = 4 \Rightarrow 1$ ————— AND
 $h_3(\text{ceg}) = 5 \Rightarrow 1$ ————— AND

mit $h_1(\text{mit}) = 8 \Rightarrow 1$ ————— AND
 $h_2(\text{mit}) = 2 \Rightarrow 1$ ————— AND
 $h_3(\text{mit}) = 1 \Rightarrow 1$ ————— AND

Anna $h_1(\text{anna}) = 8 \Rightarrow 1$ ————— AND
 $h_2(\text{anna}) = 9 \Rightarrow 0$ ————— AND
 $h_3(\text{anna}) = 5 \Rightarrow 1$ ————— AND

Estimating moments.

- A Stream consists of elements chosen from a universal set. Assume that your universal set is order. Let M_i be the number of occurrences of i^{th} element

for any i then the k^{th} order moment of the stream is $\text{Sum } (m_i)^k$

- 0th moment is the count of no. of distinct element in a stream.
- 1st moment is the sum of m_i 's must be the length of the stream.
- 2nd moment is the sum of squares of m_i 's called the surprise number which measures how uneven is the distribution of elements in the stream

Example

Stream S1: a, b, a, c, a, d, a, c

$$\begin{aligned} \text{0}^{\text{th}} \text{ order moment} &= 1^0 + 1^0 + 2^0 + 1^0 \\ &= 4 \quad (\text{no. of distinct elements}) \end{aligned}$$

$$\begin{aligned} \text{first order moment} &= 1^1 + 1^1 + 2^1 + 1^1 \\ &= 8 \quad (\text{length of the stream}) \end{aligned}$$

$$\begin{aligned} \text{second order moment} &= 1^2 + 1^2 + 2^2 + 1^2 \\ &= 22 \end{aligned}$$

Find the 1st, 0th, 2nd order moment for

10, 9, 9, 9, 9, 9, 9, 9, 9, 9

$$\begin{aligned} \text{0}^{\text{th}} \text{ order moment} &= 1^0 + 10^0 \\ &= 2 \end{aligned}$$

$$\begin{aligned} \text{1}^{\text{st}} \text{ order moment} &= 1^1 + 10^1 \\ &= 11 \end{aligned}$$

$$\begin{aligned} \text{2}^{\text{nd}} \text{ order moment} &= 1^2 + 10^2 \\ &= 101 \end{aligned}$$

Consider if the length of the stream is 11

One of the 11 elements appear 90 times
and the other 10 appears one time each. Find
the surprise number.

$$\begin{aligned}\text{surprise number} &= 90^2 + 1^2 + 1^2 + 1^2 + 1^2 + 1^2 \\ &\quad + 1^2 + 1^2 + 1^2 + 1^2 + 1^2 \\ &= 8100 + 10 \\ &= 8110\end{aligned}$$

Alon - Matias Szegedy Algorithm for
Second moments. (AMS algorithm)

- a particular element of the universal set which is referred to as a dock element
- An integer x dot value is a value of variable to ~~dock~~ determine the value of variable we choose a position in between $1 \dots n$ uniformly and at random.
Set x dot element = element and initialize x dot value = 1 as we read the stream add 1 to x dot value each time we encounter another occurrence of x dot element.
- estimate of 2nd moment
 $n(2x \text{ value} - 1)$
↓
length of the stream

Take the average for the estimation

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15
 a b c b d a c d a b d c a a b

$$\cancel{x_1 \cdot b = 4}$$

$$\cancel{x_1 \cdot d = }$$

let random positions be
 2, 8, 12

$$x_1 \cdot \text{value} = 4 \quad (\text{b}) \\ (\text{no. of times } b \text{ occurs})$$

$$x_2 \cdot \text{value} = 2 \quad (\text{d})$$

$$x_3 \cdot \text{value} = 1$$

$$x_1 \cdot \text{element} = b$$

$$x_2 \cdot \text{element} = d$$

$$x_3 \cdot \text{element} = c$$

$$\text{estimation } n [2 \times \text{value} - 1]$$

$$x_1 = 15 [2 \times 4 - 1] \\ = 105$$

$$x_2 = 15 [2 \times 2 - 1] \\ = 145$$

$$x_3 = 15 [2 \times 1 - 1] \\ = 15$$

$$\frac{105 + 145 + 15}{3} = 55$$

→ approximately equal to 2nd order moment.

Suppose we keep 3 variables x_1, x_2, x_3 , assume that at random we pick 3rd, 8th, 13th position to define these variables.

Compute 2nd moment of the stream, estimate 2nd moment of the stream and compute average of it.

a, b, c, b, d, a, c, d, a, d, b, d, c, a
a, b.

Ams algorithm

$$x_1. \text{element} = c$$

$$x_1. \text{value} = 3$$

$$x_2. \text{element} = d$$

$$x_2. \text{value} = 2$$

$$x_3. \text{element} = a$$

$$x_3. \text{value} = 2$$

estimation $n [2n \text{ value} - 1]$

$$x_1 = 15 [2 \times 3 - 1] \\ = 105$$

$$x_2 = 15 [2 \times 2 - 1] \\ = 45$$

$$x_3 = 15 [2 \times 2 - 1] \\ = 45$$

$$\text{2nd order moment} = \frac{75 + 45 + 45}{3} = 55$$

general method

$$\text{2nd order moment} = 5^2 + 4^2 + 3^2 + 3^2 \\ = 25 + 16 + 9 + 9 \\ = 59$$

Refer NQSQL

7/6/22

	m_1	m_2	m_3	m_4	m_5
u_1	3	1	1	3	1
u_2	1	2	4	1	3
u_3	3	1	1	3	?
u_4	4	3	5	4	4

	comedy	action
m_1	3	1
m_2	1	2
m_3	1	4
m_4	3	1
m_5	1	3

	comedy	action
u_1	✓	✗
u_2	✗	✓
u_3	✓	✗
u_4	✓	✓

matrix decomposition / Singular Value Decomposition

$$\begin{array}{ccc}
 m \times n & = & [m \times a] \times [a \times n] \\
 \downarrow \text{no. of users} & & \downarrow \text{rank of the matrix} \\
 \downarrow \text{no. of movies} & & (\text{hidden features/latent features})
 \end{array}$$

	m_1	m_2	m_3	m_4	m_5
u_1	3	1	1	3	1
u_2	1	2	4	1	3
u_3	3	1	1	3	1
u_4	4	3	5	4	4

$m \times n$
 4×5

Model Based recommendation System

Advantages of SVD

- storage is reduced

$(4 \times 5 = 20 \text{ entries})$
 By SVD ($5 \times 2 + 4 \times 2 = 18 \text{ entries}$)

key value database
 e.g. product details
 key userid

user id
 Username
 Orderid
 Product name
 Product id

9/6/22

Document store

- eg MongoDB
- Can be stored as PDF, word, JSON, formats

FIND.id

Tabular store

- eg HBASE

HBASE does not support join operation.

HIVE

- HQL is used to write map-reduce function
- time consuming, thus not suitable for real time applications.

Module 10

10/6/22

Different types of visualization

- * basic
- * advanced

histogram (data distribution)

box plot (provides 5 number summary)

bar plot (compare categorical data)

scatterplot (find relationship between two features)

lineplot (suitable for timeseries data eg stockmarket data)
analyse the trend

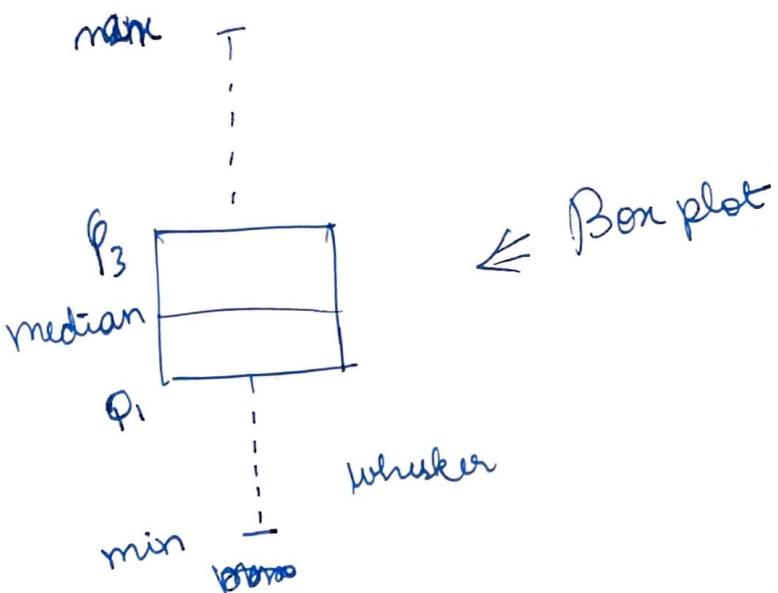
piechart (find the contribution of

GGplot - grammar of graphics.

4 arguments (dataset, aesthetics, geometries, facets)

mapping of
x, y
coordinates

per
bar
scatter
etc.



attributes such as
aesthetics map to values of n labels, y labels, shape,
color

geometries - show the data is displayed.
type of data visualisation

facets - display the subset of data using
columns and rows

Titanic dataset

~~Titanic~~

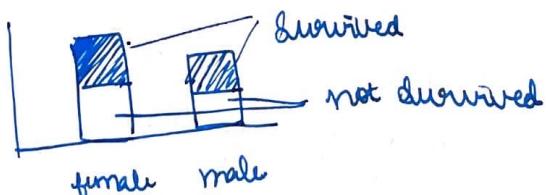
1. survival rate
2. survival by gender
3. survival rate by passenger class
4. Age distribution
5. survival rate by age

bar graph \leftarrow ggplot(titanic ~ as(x = survived)) + geom_bar()

prop.table(table(titanic \$survived))

gives the percentage of each unique
value in survived

ggplot(titanic ~ as(x = psex, fill = survived)) + geom_bar()



Review of data analytics using R

as factor to convert to categorical variable

exploratory data analysis

- different visualization methods
- ggplot

summary()

str() - structure of dataset

head()

lm - linear model

Graph memory

graph

- used for social media / maps / friendship network / page rank algorithm.

influence mapping

find centrality .

Types of graph

embedded micronetworks

communication model

collaborative communities

influence modeling

distance modeling

important metrics

degree - no. of in connections and out connections,

betweenness -

Centrality - more influential person in the network.

Closeness

Fair Fairness

Density