

SQL MASTERCLASS

BARCELONA DATA INSTITUTE

Iago Novoa
Head of BI at letgo

WRITE IN A PAPER

FAVOURITE PLACE ON EARTH

SQL MasterClass
Barcelona Data Institute

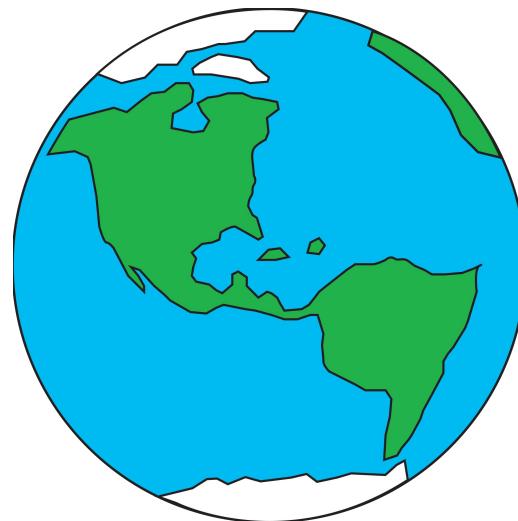
1. INTRODUCTION TO DATABASES
2. OUR DATA WORLD
3. AMAZON REDSHIFT VS GOOGLE BIGQUERY
4. HOW TO CREATE A DATABASE
5. BASIC SQL
6. JOIN-ING TABLES
7. WINDOW FUNCTIONS
8. HOW TO EXPORT AND IMPORT DATA IN REDSHIFT
9. QA

SQL MasterClass
Barcelona Data Institute

1. INTRODUCTION TO DATABASES
2. OUR DATA WORLD
3. AMAZON REDSHIFT VS GOOGLE BIGQUERY
4. HOW TO CREATE A DATABASE
5. BASIC SQL
6. JOIN-ING TABLES
7. WINDOW FUNCTIONS
8. HOW TO EXPORT AND IMPORT DATA IN REDSHIFT
9. QA

SQL MasterClass
Barcelona Data Institute

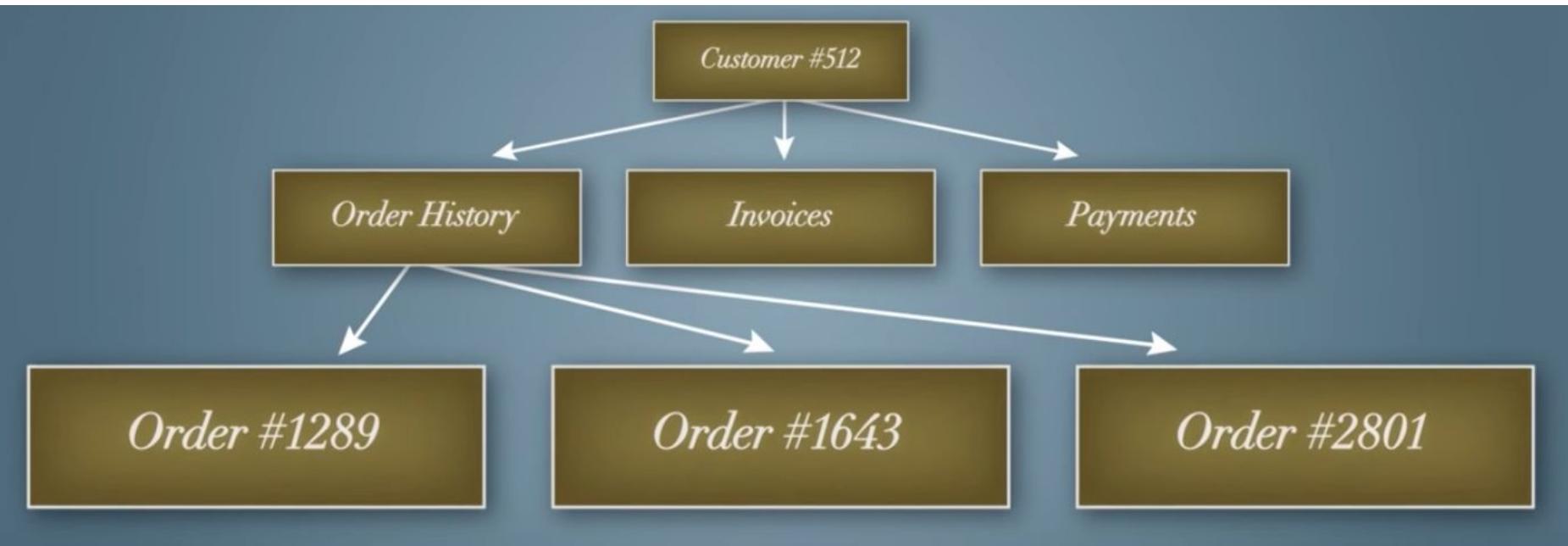
THOSE OF US CURRENTLY ALIVE REPRESENT ~7% OF
TOTAL NUMBER OF HUMANS WHO HAVE EVER LIVED



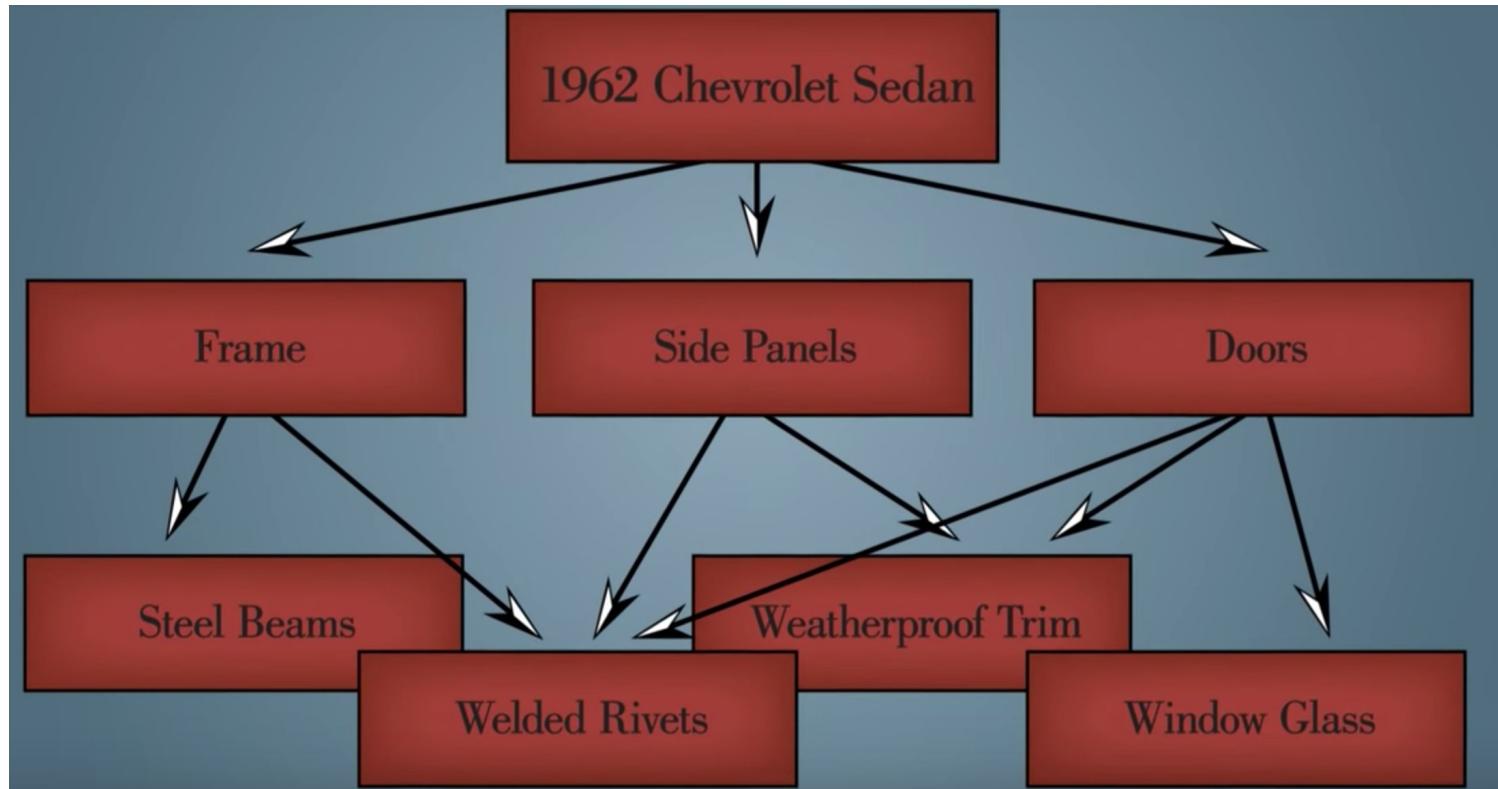
HUMAN BEINGS BEGAN TO STORE INFORMATION VERY LONG AGO



COMPUTERIZED DBS STARTED IN 1960S WITH HIERARCHICAL MODEL ...



... AND NETWORK MODELS



IN THE 1970S CODD PROPOSED THE USE OF RELATIONAL MODEL

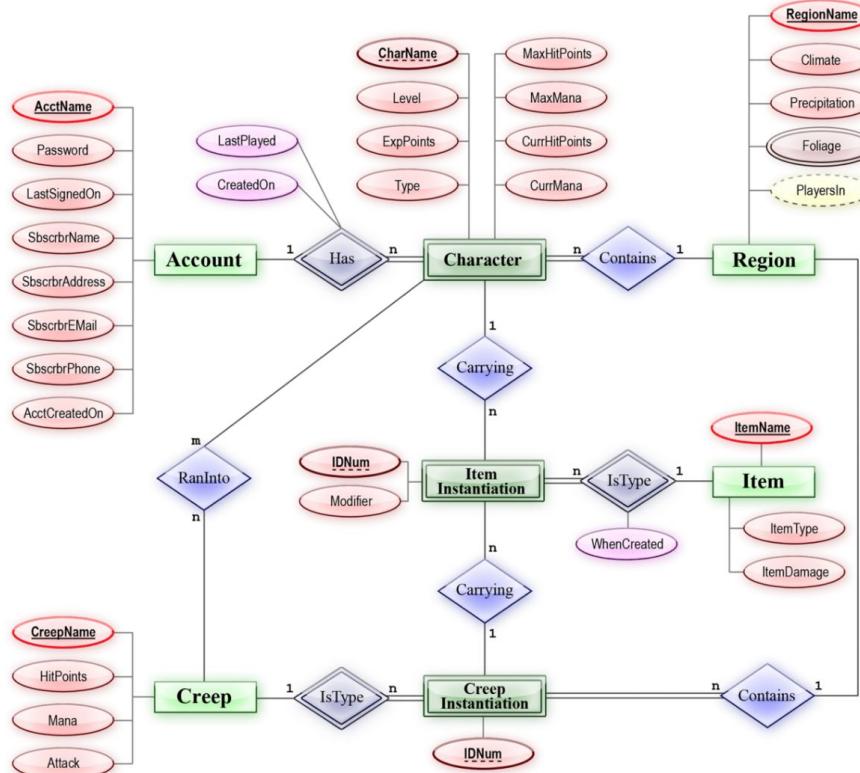
Relational Model

Movies		
Movie Name	Studio ID	Actor ID
Towering Inferno	20F	31
The Godfather	PP	22
Chinatown	PP	14
Straw Dogs	ABC	57

Actors		
Actor ID	Actor Name	Born
57	Dustin Hoffman	1937
42	Al Pacino	1940
14	Jack Nicholson	1937
22	Diane Keaton	1946

Studios			
Studio ID	Studio Name	Founded	Headquarters
PP	Paramount Pictures	1912	Hollywood
UA	United Artists	1919	Los Angeles
20F	20th Century Fox	1935	Century City

AND CHEN PROPOSED ENTITY-RELATIONSHIP (ER) MODEL

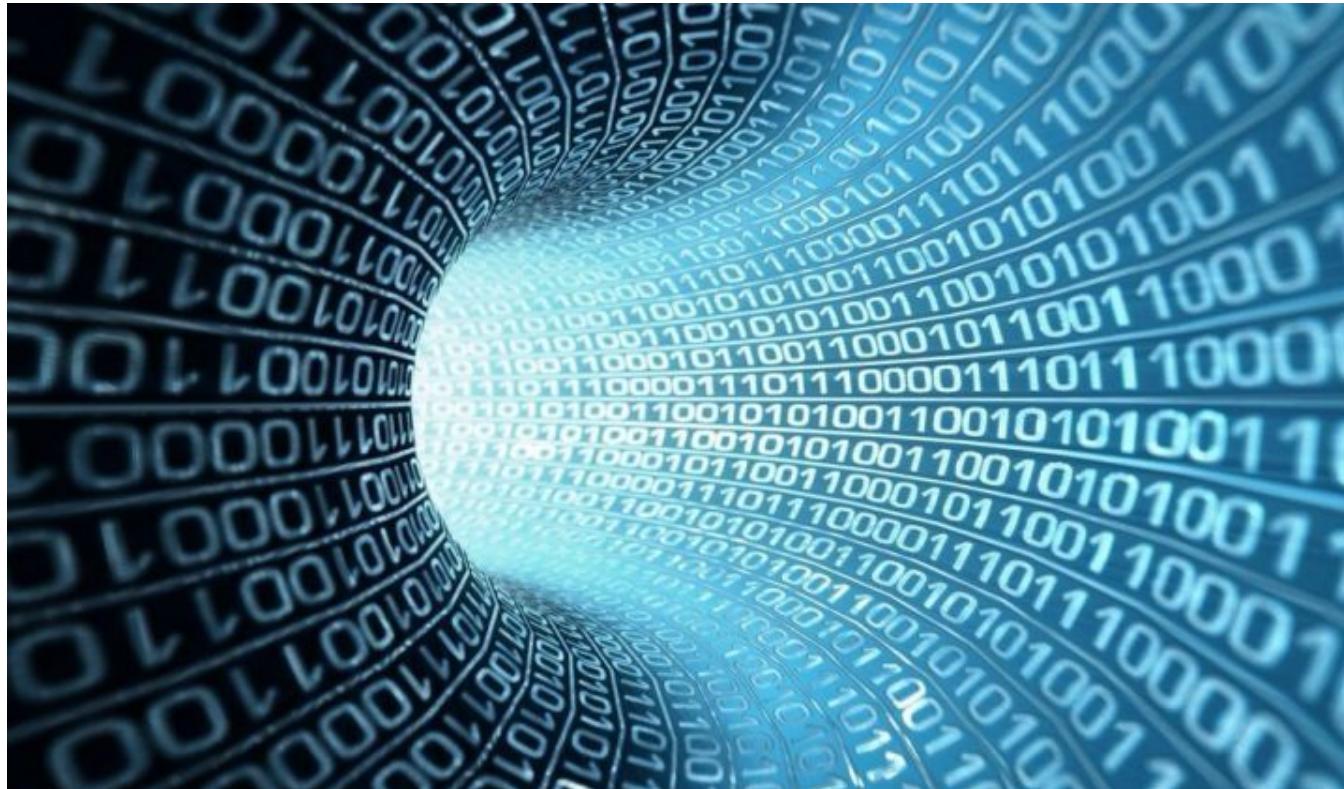


IN 1980S SQL BECAME THE STANDARD QUERY LANGUAGE

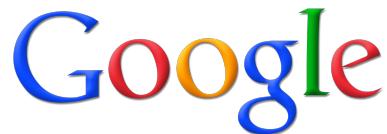
```
CREATE TABLE test ( a INTEGER, b TEXT, c T  
SELECT * FROM  
    [.....]  
    INTO test  
    INTO test (c ) VALUES  
    INTO test (b, c ) SEL
```



INTERNET LED TO EXPONENTIAL GROWTH OF THE DATABASE INDUSTRY

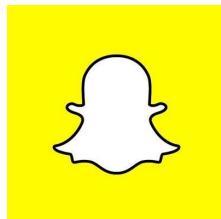


2.5 QUINTILLION BYTES OF DATA CREATED EACH DAY



Google now processes more than
3.5 billion searches per day

Over the last two years alone **90** percent of the data in the world was generated



Snapchat users share **527,760** photos every minute

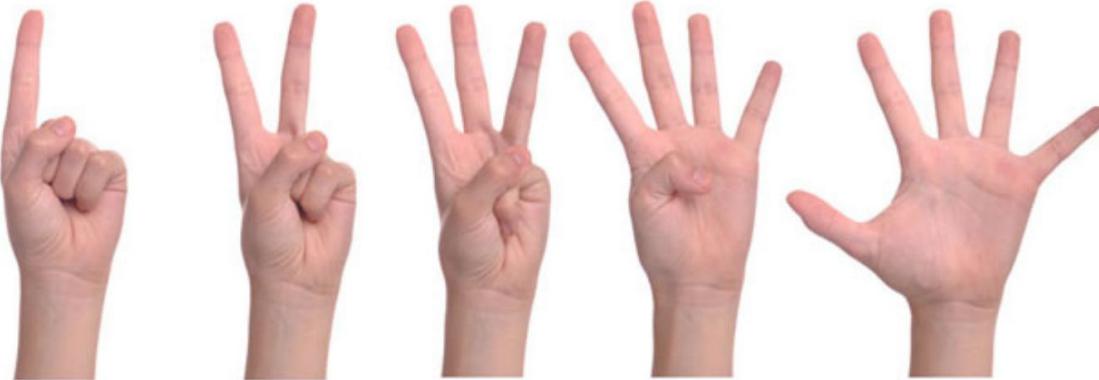


YouTube

Users watch
248,796,600 YouTube videos per hour

DATABASE MANAGEMENT SYSTEM (DBMS) IS A SOFTWARE TO DEFINE, BUILD, MAINTAIN, ... A DATABASE





I don't
understand
at all

I need to
go over
this again

I think I got
it, but am
not
completely
comfortable

I got it

I can
explain it
to someone
else

1. INTRODUCTION TO DATABASES
2. OUR DATA WORLD
3. AMAZON REDSHIFT VS GOOGLE BIGQUERY
4. HOW TO CREATE A DATABASE
5. BASIC SQL
6. JOIN-ING TABLES
7. WINDOW FUNCTIONS
8. HOW TO EXPORT AND IMPORT DATA IN REDSHIFT
9. QA

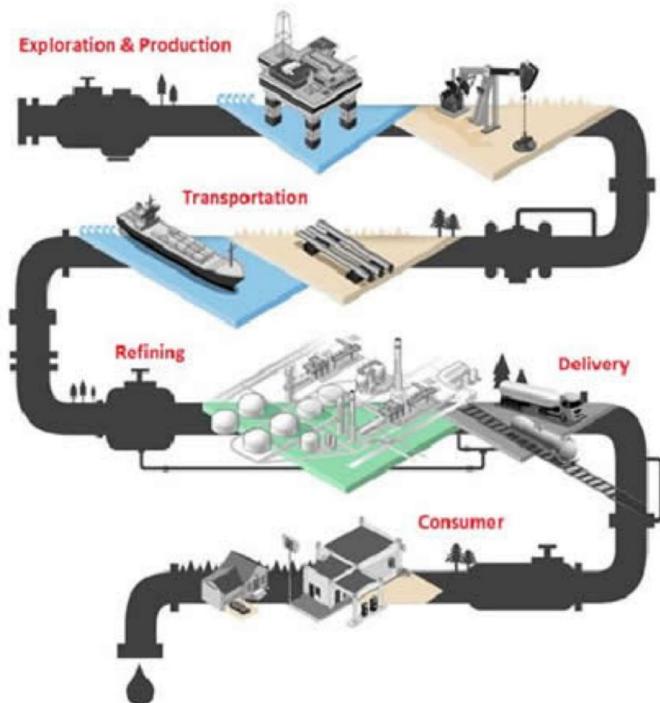
SQL MasterClass
Barcelona Data Institute



DATA IS THE MAIN DRIVER OF MANY DEALS NOWADAYS

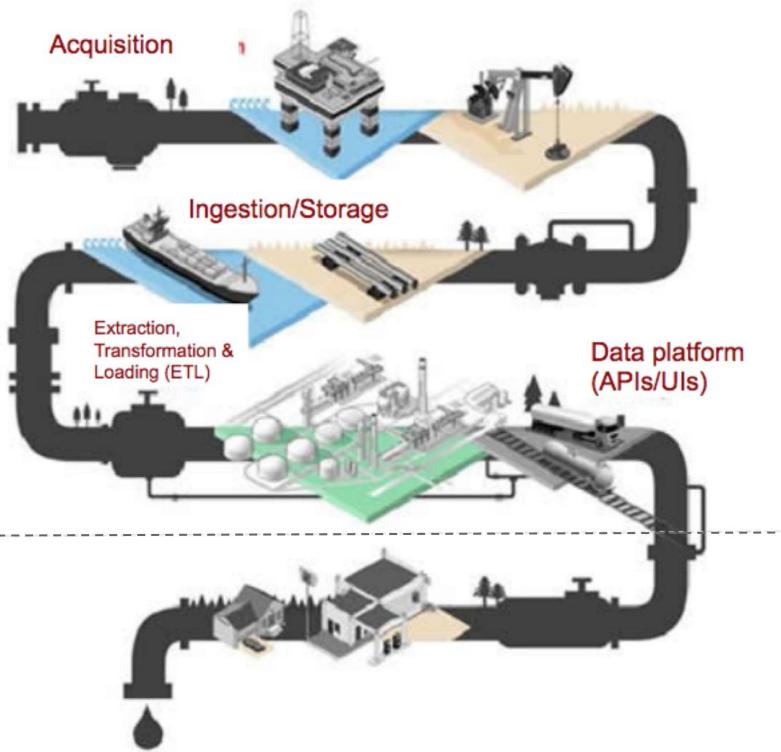
	Target company (date)	Value of deal (\$USbn)	Business
facebook	INSTAGRAM (2012)	1	1.0 PHOTO SHARING
	WHATSAPP (2014)	?	TEXT/PHOTO MESSAGING
Alphabet	WAZE (2013)	1	1.2 MAPPING AND NAVIGATION
IBM	THE WEATHER COMPANY (2015)	1	2.0 METEOROLOGY
	TRUVEN HEALTH ANALYTICS (2016)	1	2.6 HEALTH CARE
intel	MOBILEYE (2017)	15.3	SELF-DRIVING CARS
Microsoft	SWIFTKEY (2016)	0.25	KEYBOARD/AI
	LINKEDIN (2016)	?	BUSINESS NETWORKING
ORACLE	BLUEKAI (2014)	0.4	CLOUD DATA PLATFORM
	DATACOGIX (2014)	1.0	MARKETING

THE OIL GENERATION



THE DATA GENERATION

DATA ENGINEER



DATA ANALYST / DATA SCIENTIST

BUILDING WELL-BALANCED DATA PRODUCT TEAMS



"VW BEETLE JET" TEAM

- 3 DATA ANALYSTS / DATA SCIENTISTS
- 2 DATA ENGINEERS
- 1 PRODUCT MANAGER

- NOISY
- NOT FUEL EFFICIENT
- CAN EXPLODE

BUILDING WELL-BALANCED DATA PRODUCT TEAMS



"PRIUS" TEAM

- 1 DATA ANALYST / DATA SCIENTIST
- 1 DATA ENGINEER
- 1 BACKEND ENGINEER
- 1 PRODUCT MANAGER

- Nothing fancy
- Very efficient
- Will definitely not explode



I don't
understand
at all

I need to
go over
this again

I think I got
it, but am
not
completely
comfortable

I got it

I can
explain it
to someone
else

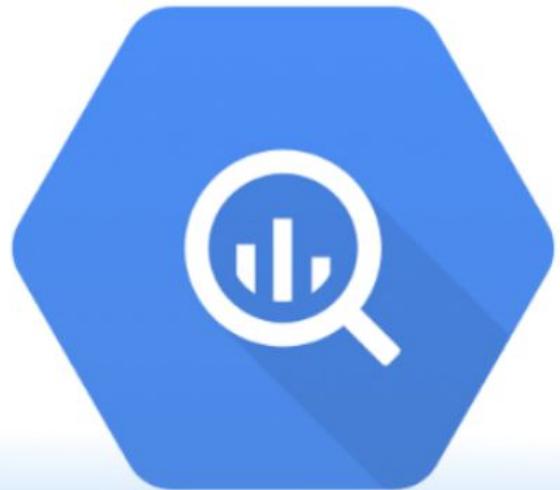
1. INTRODUCTION TO DATABASES
2. OUR DATA WORLD
3. **AMAZON REDSHIFT VS GOOGLE BIGQUERY**
4. HOW TO CREATE A DATABASE
5. BASIC SQL
6. JOIN-ING TABLES
7. WINDOW FUNCTIONS
8. HOW TO EXPORT AND IMPORT DATA IN REDSHIFT
9. QA

SQL MasterClass
Barcelona Data Institute

MCDONALD'S OR BURGER KING? NIKE OR ADIDAS?



VS



AMAZON REDSHIFT VS GOOGLE BIGQUERY

Amazon Redshift	Google BigQuery
Fork of PostgreSQL	Was built from scratch
Relational	Nested data structures are first class citizen
SQL	NoSQL (but now it also speaks SQL)
Analysts are more familiar with it	Developers love it

LOADING DATA ...

	Amazon Redshift	Google BigQuery
Bulk data:	?	Google Cloud Storage
Streaming data:	Kinesis	Streaming Inserts
Other sources:		Google Analytics Premium
Data Serializations:	CSV, Avro, JSON	CSV, Avro, JSON

DATA MODELING

	Amazon Redshift	Google BigQuery
Data organization:	Tables, Schemas	Tables, Datasets
Datatypes:	Subset of SQL datatypes	Custom Datatypes with mappings
Nested datatypes:	Not supported	Natively supported
Table manipulation:	Supported with SQL.	?
Support for updates and deletes:	Natively through SQL with no additional costs.	Limited and expensive.



DATA CONSISTENCY

	Amazon Redshift	Google BigQuery
Transactions:	Yes	No
Deduplication (batch)	Easy to do	?
Deduplication (Streaming)	No	Yes with time window + ID
Delivery Semantics (Streaming)	At least once (Kinesis)	At least once with best effort deduplication.

OTHER THINGS TO TAKE INTO ACCOUNT

- CLUSTER MANAGEMENT
- CONNECTIVITY
- AUTHENTICATION
- RESOURCE ALLOCATION
- DATABASE OPTIMIZATION
- COST
- ECOSYSTEM
- ...







I don't
understand
at all

I need to
go over
this again

I think I got
it, but am
not
completely
comfortable

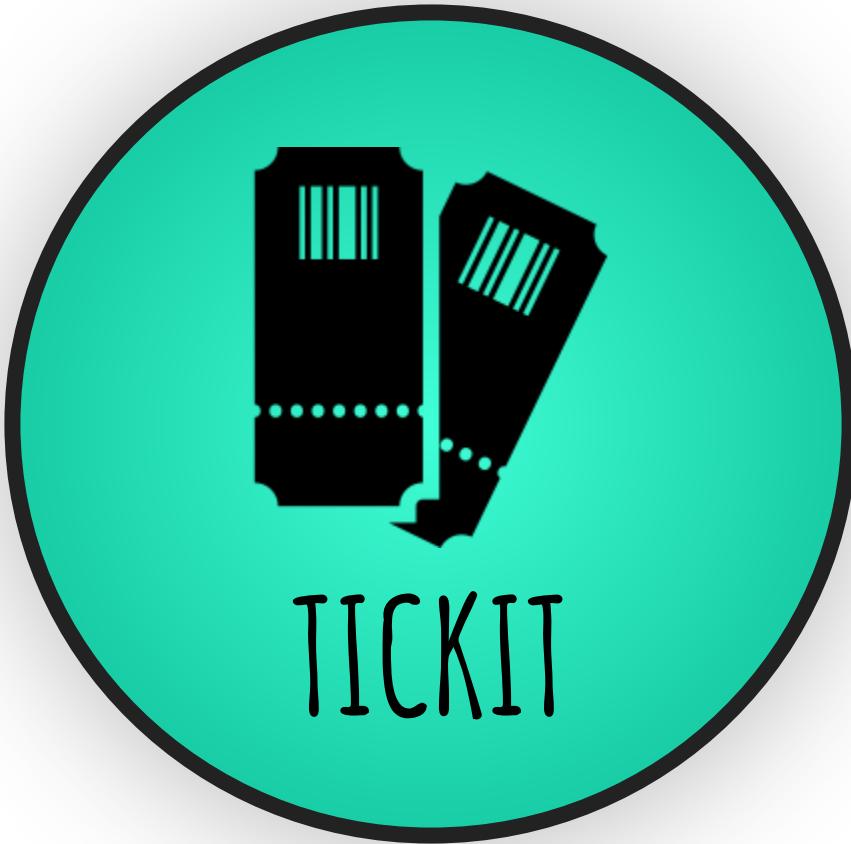
I got it

I can
explain it
to someone
else

1. INTRODUCTION TO DATABASES
2. OUR DATA WORLD
3. AMAZON REDSHIFT VS GOOGLE BIGQUERY
4. HOW TO CREATE A DATABASE
5. BASIC SQL
6. JOIN-ING TABLES
7. WINDOW FUNCTIONS
8. HOW TO EXPORT AND IMPORT DATA IN REDSHIFT
9. QA

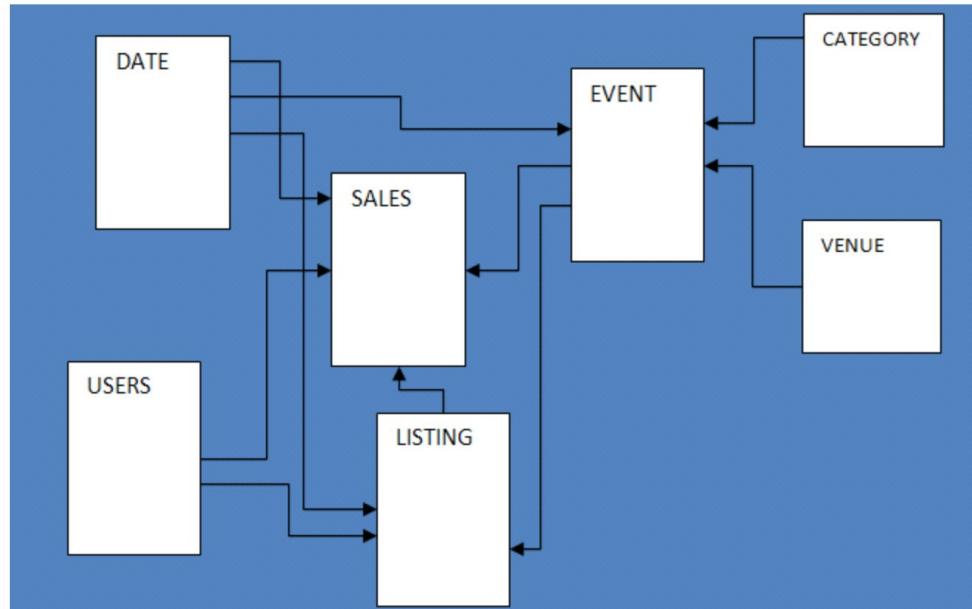
SQL MasterClass
Barcelona Data Institute





1. DESIGN OUR DATABASE

This sample database application helps analysts track sales activity for the fictional TICKIT web site, where users buy and sell tickets online for sporting events, shows, and concerts. In particular, analysts can identify ticket movement over time, success rates for sellers, and the best-selling events, venues, and seasons. Analysts can use this information to provide incentives to buyers and sellers who frequent the site, to attract new users, and to drive advertising and promotions.



2. DEFINE TABLES

CATEGORY Table

Column Name	Data Type	Description
CATID	SMALLINT	Primary key, a unique ID value for each row. Each row represents a specific type of event for which tickets are bought and sold.
CATGROUP	VARCHAR(10)	Descriptive name for a group of events, such as Shows and Sports .
CATNAME	VARCHAR(10)	Short descriptive name for a type of event within a group, such as Opera and Musicals .
CATDESC	VARCHAR(30)	Longer descriptive name for the type of event, such as Musical theatre .

SALES Table

Column Name	Data Type	Description
SALESID	INTEGER	Primary key, a unique ID value for each row. Each row represents a sale of one or more tickets for a specific event, as offered in a specific listing.
LISTID	INTEGER	Foreign-key reference to the LISTING table.
SELLERID	INTEGER	Foreign-key reference to the USERS table (the user who sold the tickets).
BUYERID	INTEGER	Foreign-key reference to the USERS table (the user who bought the tickets).
EVENTID	INTEGER	Foreign-key reference to the EVENT table.
DATEID	SMALLINT	Foreign-key reference to the DATE table.
QTY SOLD	SMALLINT	The number of tickets that were sold, from 1 to 8. (A maximum of 8 tickets can be sold in a single transaction.)
PRICEPAID	DECIMAL(8,2)	The total price paid for the tickets, such as 75.00 or 488.00 . The individual price of a ticket is PRICEPAID/QTY SOLD.
COMMISSION	DECIMAL(8,2)	The 15% commission that the business collects from the sale, such as 11.25 or 73.20 . The seller receives 85% of the PRICEPAID value.
SALETIME	TIMESTAMP	The full date and time when the sale was completed, such as 2008-05-24 06:21:47 .

DATE Table

Column Name	Data Type	Description
DATEID	SMALLINT	Primary key, a unique ID value for each row. Each row represents a day in the calendar year.
CALDATE	DATE	Calendar date, such as 2008-06-24 .
DAY	CHAR(3)	Day of week (short form), such as SA .
WEEK	SMALLINT	Week number, such as 26 .
MONTH	CHAR(5)	Month name (short form), such as JUN .
QTR	CHAR(5)	Quarter number (1 through 4).
YEAR	SMALLINT	Four-digit year (2008).
HOLIDAY	BOOLEAN	Flag that denotes whether the day is a public holiday (U.S.).

EVENT Table

Column Name	Data Type	Description
EVENTID	INTEGER	Primary key, a unique ID value for each row. Each row represents a separate event that takes place at a specific venue at a specific time.
VENUEID	SMALLINT	Foreign-key reference to the VENUE table.
CATID	SMALLINT	Foreign-key reference to the CATEGORY table.
DATEID	SMALLINT	Foreign-key reference to the DATE table.
EVENTNAME	VARCHAR(200)	Name of the event, such as Hamlet or La Traviata .
STARTTIME	TIMESTAMP	Full date and start time for the event, such as 2008-10-10 19:30:00 .

VENUE Table

Column Name	Data Type	Description
VENUEID	SMALLINT	Primary key, a unique ID value for each row. Each row represents a specific venue where events take place.
VENUENAME	VARCHAR(100)	Exact name of the venue, such as Cleveland Browns Stadium .
VENUECITY	VARCHAR(30)	City name, such as Cleveland .
VENUESTATE	CHAR(2)	Two-letter state or province abbreviation (United States and Canada), such as OH .
VENUESEATS	INTEGER	Maximum number of seats available at the venue, if known, such as 73200 . For demonstration purposes, this column contains some null values and zeroes.

LISTING Table

Column Name	Data Type	Description
LISTID	INTEGER	Primary key, a unique ID value for each row. Each row represents a listing of a batch of tickets for a specific event.
SELLERID	INTEGER	Foreign-key reference to the USERS table, identifying the user who is selling the tickets.
EVENTID	INTEGER	Foreign-key reference to the EVENT table.
DATEID	SMALLINT	Foreign-key reference to the DATE table.
NUMTICKETS	SMALLINT	The number of tickets available for sale, such as 2 or 20 .
PRICEPERTICKET	DECIMAL(8,2)	The fixed price of an individual ticket, such as 27.00 or 206.00 .
TOTALPRICE	DECIMAL(8,2)	The total price for this listing (NUMTICKETS*PRICEPERTICKET).
LISTTIME	TIMESTAMP	The full date and time when the listing was posted, such as 2008-03-18 07:19:35 .

USERS Table

Column Name	Data Type	Description
USERID	INTEGER	Primary key, a unique ID value for each row. Each row represents a registered user (a buyer or seller or both) who has listed or bought tickets for at least one event.
USERNAME	CHAR(8)	An 8-character alphanumeric username, such as POL0BLJZ .
FIRSTNAME	VARCHAR(30)	The user's first name, such as Victor .
LASTNAME	VARCHAR(30)	The user's last name, such as Hernandez .
CITY	VARCHAR(30)	The user's home city, such as Kaperville .
STATE	CHAR(2)	The user's home state, such as GA .
EMAIL	VARCHAR(100)	The user's email address; this column contains random Latin values, such as turpiseaccumanlaureet.org .
PHONE	CHAR(14)	The user's 14-character phone number, such as (818) 765-4255 .
LIKESPORTS	BOOLEAN	A series of 10 different columns that identify the user's likes and dislikes with true and false values.

3. CREATE TABLES



```
create table users(
    userid integer not null distkey sortkey,
    username char(8),
    firstname varchar(30),
    lastname varchar(30),
    city varchar(30),
    state char(2),
    email varchar(100),
    phone char(14),
    likesports boolean,
    liketheatre boolean,
    likeconcerts boolean,
    likejazz boolean,
    likeclassical boolean,
    likeopera boolean,
    likerock boolean,
    likevegas boolean,
    likebroadway boolean,
    likemusicals boolean);

create table venue(
    venueid smallint not null distkey sortkey,
    venuename varchar(100),
    venuecity varchar(30),
    venuestate char(2),
```

4. LOAD DATA



```
copy users from 's3://awssampledbuswest2/ticket/allusers_pipe.txt'
credentials 'aws_iam_role=<iam-role-arn>'
delimiter '||' region 'us-west-2';

copy venue from 's3://awssampledbuswest2/ticket/venue_pipe.txt'
credentials 'aws_iam_role=<iam-role-arn>'
delimiter '||' region 'us-west-2';

copy category from 's3://awssampledbuswest2/ticket/category_pipe.txt'
credentials 'aws_iam_role=<iam-role-arn>'
delimiter '||' region 'us-west-2';

copy date from 's3://awssampledbuswest2/ticket/date2008_pipe.txt'
credentials 'aws_iam_role=<iam-role-arn>'
delimiter '||' region 'us-west-2';

copy event from 's3://awssampledbuswest2/ticket/allevents_pipe.txt'
credentials 'aws_iam_role=<iam-role-arn>'
delimiter '||' timeformat 'YYYY-MM-DD HH:MI:SS' region 'us-west-2';

copy listing from 's3://awssampledbuswest2/ticket/listings_pipe.txt'
credentials 'aws_iam_role=<iam-role-arn>'
delimiter '||' region 'us-west-2';

copy sales from 's3://awssampledbuswest2/ticket/sales_tab.txt'
credentials 'aws_iam_role=<iam-role-arn>'
delimiter '\t' timeformat 'MM/DD/YYYY HH:MI:SS' region 'us-west-2';
```

NOW WE ARE READY TO START QUERYING...

```
select * from sql_masterclass.sales;
```



```
select * from sql_masterclass.sales limit 10;
```





I don't
understand
at all

I need to
go over
this again

I think I got
it, but am
not
completely
comfortable

I got it

I can
explain it
to someone
else

1. INTRODUCTION TO DATABASES
2. OUR DATA WORLD
3. AMAZON REDSHIFT VS GOOGLE BIGQUERY
4. HOW TO CREATE A DATABASE
5. **BASIC SQL**
6. JOIN-ING TABLES
7. WINDOW FUNCTIONS
8. HOW TO EXPORT AND IMPORT DATA IN REDSHIFT
9. QA

SQL MasterClass
Barcelona Data Institute

SQL IS A KEY SKILL FOR ANY DATA ANALYST



SELECT ...

GET THE NAME OF 5 EVENTS

```
select
    eventname
from |
    sql_masterclass.event
limit 5;
```

Result

	eventname
1	Mamma Mia!
2	Return To Forever
3	Hannah Montana
4	Smashing Pumpkins
5	K.D. Lang

BASIC ELEMENTS

GET ALL INFORMATION FROM CATEGORIES WITHIN SPORTS GROUP

```
select * from sql_masterclass.category where catgroup='Sports';
```

Result

	catid	catgroup	catname	catdesc
1	1	Sports	MLB	Major League Baseball
2	2	Sports	NHL	National Hockey League
3	3	Sports	NFL	National Football League
4	4	Sports	NBA	National Basketball Association
5	5	Sports	MLS	Major League Soccer

EXPRESSION LISTS

GET NAME OF ALL THE VENUES IN FLORIDA AND OHIO

```
select
    venuename
from
    sql_masterclass.venue
where
    venuestate IN ('FL', 'OH');
```

Result

	venuename
1	Quicken Loans Arena
2	American Airlines Arer
3	Progressive Field
4	BankAtlantic Center
5	Nationwide Arena
6	Dolphin Stadium
7	Paul Brown Stadium
8	Raymond James Stadi
9	Tropicana Field
10	E.J. Nutter Center
11	Great American Ball Pa
12	Columbus Crew Stadi
13	Amway Arena
14	St. Pete Times Forum
15	Cleveland Browns Stac
16	Jacksonville Municipal

MATCHING CONDITIONS

GET FULL NAME (FIRST AND LAST NAME) OF 10 USERS WITH EMAILS THAT END BY 'NUNC.EDU'

```
select
    firstname || ' ' || lastname as full_name,
    email
from
    sql_masterclass.users
where
    email like '%nunc.edu'
limit 10;
```

Result

select firstname || '' || lastname as | Enter a SQL expression

	full_name	email
1	Lana Ross	mi.lorem@nunc.edu
2	Anika Camacho	nibh.Aliquam@Fusce
3	Forrest Dunn	augue.scelerisque.m
4	Wesley Bowen	aliquam.adipiscing@i
5	Alma Bean	In@consectetueripsu
6	Jescie Pickett	lobortis@Sednunc.ed
7	Yoko Burgess	vitae.velit@maurisanc
8	Xavier Mccray	velit.Pellentesque@ni
9	Medge Witt	Sed.nunc.est@eratSe
10	Rylee Cummings	rutrum.magna.Cras@

COMPARISON CONDITIONS

GET FULL NAME OF THE VENUE AND NUMBER OF SEATS FOR VENUES WITH MORE THAN 1.000 SEATS
ORDERING BY NUMBER OF SEATS (HIGHEST FIRST)

```
select
    venueusername,
    venueseats
from
    sql_masterclass.venue
where
    venueseats > 10000
order by venueseats desc;
```

Result

	venueusername	venueseats
1	FedExField	91,704
2	New York Giants Stadium	80,242
3	Arrowhead Stadium	79,451
4	INVESCO Field	76,125
5	Dolphin Stadium	74,916
6	Ralph Wilson Stadium	73,967
7	Jacksonville Municipal Stadium	73,800
8	Bank of America Stadium	73,298
9	Cleveland Browns Stadium	73,200
10	Lambeau Field	72,922
11	Reliant Stadium	72,000
12	Louisiana Superdome	72,000

LOGICAL CONDITIONS

GET USERNAME OF 10 USERS THAT LIKE MUSICAL AND THEY ALSO LIKE JAZZ OR OPERA

```
select
    username,
    likemusicals,
    likejazz,
    likeopera
from
    sql_masterclass.users
where
    likemusicals=true and (likejazz=true or likeopera=true)
limit 10;
```

Result

	username	likemusicals	likejazz	likeopera
1	NDQ15VBM	true	true	true
2	TNP95CNK	true	[NULL]	true
3	UKK70QBF	true	[NULL]	true
4	MWK08RDN	true	[NULL]	true
5	UJE76BTX	true	[NULL]	true
6	MVZ73SEM	true	true	false
7	BFL14DIV	true	false	true
8	NBV80XDJ	true	[NULL]	true
9	RJK72ITL	true	true	[NULL]
10	QSF62BKU	true	true	true

AGGREGATE FUNCTIONS

GET AVERAGE AND MEDIAN OF PRICES IN SALES

```
select  
    avg(pricepaid),  
    median(pricepaid)  
from  
    sql_masterclass.sales;
```

Result

select avg(pricepaid), median(pricepaid)

	avg	median
1	642.28	386

GROUP BY AND ORDER BY

GET NUMBER OF USERS PER STATE

```
select
    state,
    count(1)
from
    sql_masterclass.users
group by 1;
```

Result

	state	count
1	SK	1,835
2	WV	481
3	MB	1,916
4	BC	1,958
5	NU	1,883
6	MO	525
7	NL	1,919
8	CO	493

GET NUMBER OF USERS PER STATE ORDERING BY BIGGEST FIRST

```
select
    state,
    count(1)
from
    sql_masterclass.users
group by 1
order by 2 desc;
```

Result

	state	count
1	NT	1,998
2	NB	1,960
3	BC	1,958
4	QC	1,929
5	NL	1,919
6	YT	1,919
7	MB	1,916
8	PE	1,906

HAVING

GET NUMBER OF USERS PER STATE ORDERING BY BIGGEST FIRST ONLY FOR THOSE STATES WITH MORE THAN 1,950 USERS

```
select
    state,
    count
from
    (select
        state,
        count(1)
    from
        sql_masterclass.users
    group by 1)
where count>1950
order by 2 desc;
```

Result

	state	count
1	NT	1,998
2	NB	1,960
3	BC	1,958

```
select
    state,
    count(1)
from
    sql_masterclass.users
group by 1
having count(1)>1950
order by 2 desc;
```

Result

	state	count
1	NT	1,998
2	NB	1,960
3	BC	1,958

CONDITIONAL FUNCTIONS

GET NUMBER OF USERS IN PENNSYLVANIA, WASHINGTON AND THE REST OF THE COUNTRY

```
select
  case
    when state='PA' then 'Pennsylvania'
    when state='WA' then 'Washington'
    else 'Rest of the Country'
  end as location,
  count(distinct userid) as users
from
  sql_masterclass.users
group by 1;
```

Result

select case when state='PA' then 'P' | Enter a SQL expression

	location	users
1	Pennsylvania	497
2	Rest of the Country	49,004
3	Washington	489

HOMEWORK 1

TASK:

- GET FIRST NAME, LAST NAME AND INITIALS (FIRST LETTER OF THE FIRST NAME AND FIRST LETTER OF LAST NAME) FROM 10 USERS
- SAVE THE QUERY AS HOMEWORK1.SQL IN GITHUB IN SQL_MASTERCLASS FOLDER

TIPS:

- USE STRING FUNCTIONS

	firstname	lastname	initials
1	Lars	Ratliff	LR
2	Colton	Roy	CR
3	Bruce	Beck	BB
4	Henry	Cochran	HC
5	Cody	Moss	CM
6	Ralph	Bird	RB
7	Wing	Jennings	WJ
8	Guy	Cochran	GC
9	Emerald	Chan	EC
10	Gloria	Rodriguez	GR

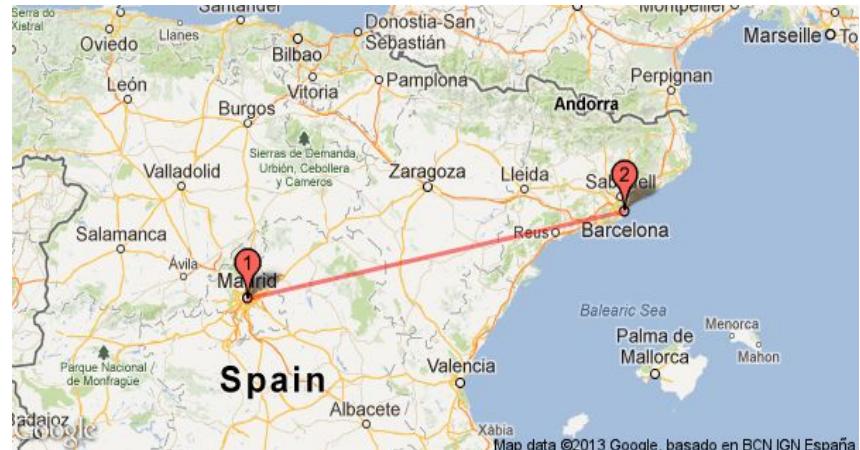
HOMEWORK 2

TASK:

- CALCULATE DISTANCE IN KM BETWEEN MADRID (41.3851 N, 2.1734 E) AND BARCELONA (40.4168 N, 3.7038 W)
- SAVE THE QUERY AS HOMEWORK2.SQL IN GITHUB IN SQL_MASTERCLASS FOLDER

TIPS:

- LOOK FOR A MATH FORMULA TO CALCULATE DISTANCE BETWEEN 2 POINTS





I don't
understand
at all

I need to
go over
this again

I think I got
it, but am
not
completely
comfortable

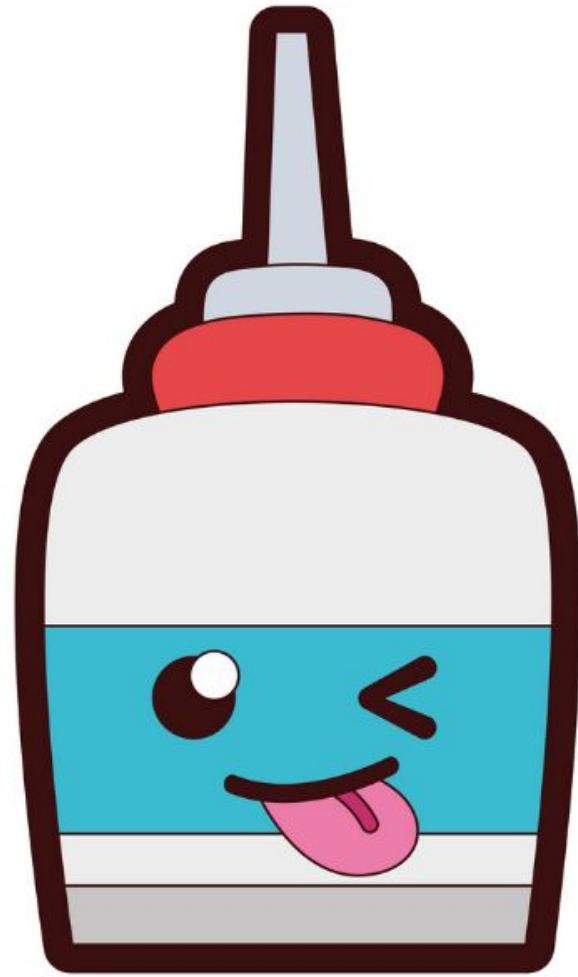
I got it

I can
explain it
to someone
else

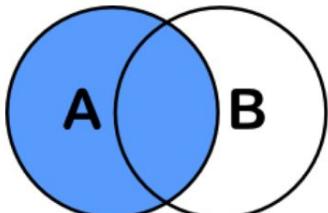


1. INTRODUCTION TO DATABASES
2. OUR DATA WORLD
3. AMAZON REDSHIFT VS GOOGLE BIGQUERY
4. HOW TO CREATE A DATABASE
5. BASIC SQL
6. JOIN-ING TABLES
7. WINDOW FUNCTIONS
8. HOW TO EXPORT AND IMPORT DATA IN REDSHIFT
9. QA

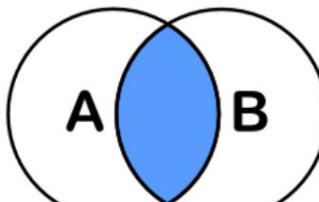
SQL MasterClass
Barcelona Data Institute



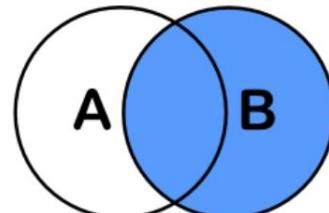
CHEATSHEET
SQL
JOINS



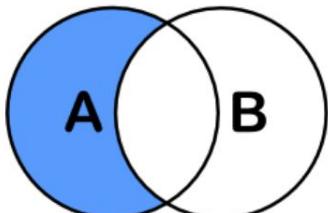
```
SELECT <auswahl>
FROM tabelleA A
LEFT JOIN tabelleB B
ON A.key = B.key
```



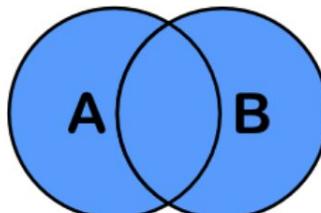
```
SELECT <auswahl>
FROM tabelleA A
INNER JOIN tabelleB B
ON A.key = B.key
```



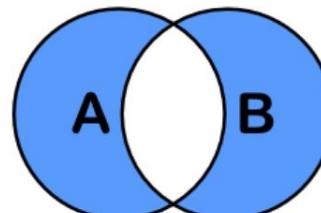
```
SELECT <auswahl>
FROM tabelleA A
RIGHT JOIN tabelleB B
ON A.key = B.key
```



```
SELECT <auswahl>
FROM tabelleA A
LEFT JOIN tabelleB B
ON A.key = B.key
WHERE B.key IS NULL
```



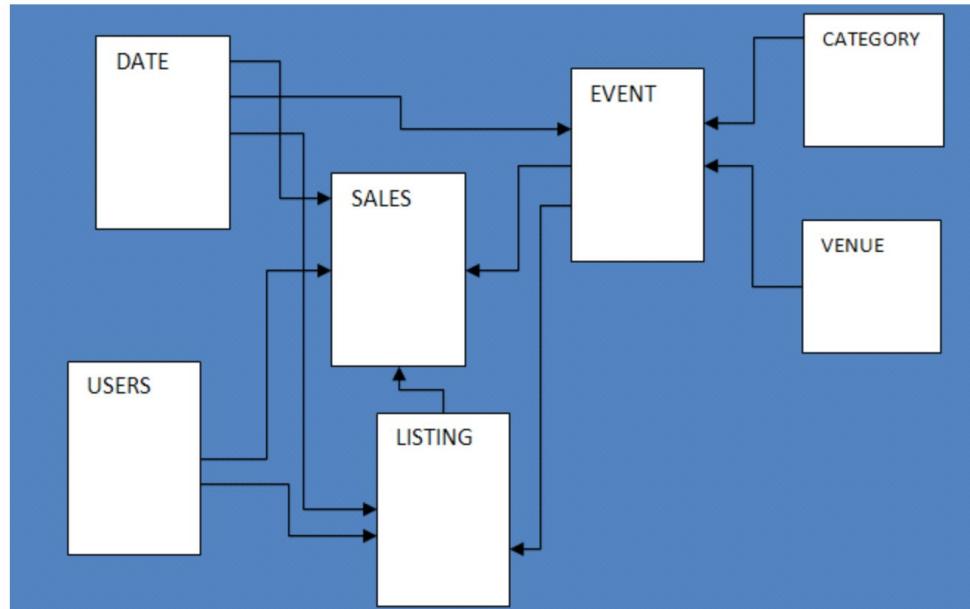
```
SELECT <auswahl>
FROM tabelleA A
FULL OUTER JOIN tabelleB B
ON A.key = B.key
```



```
SELECT <auswahl>
FROM tabelleA A
FULL OUTER JOIN tabelleB B
ON A.key = B.key
WHERE A.key IS NULL
OR B.key IS NULL
```

OUR DATABASE

This sample database application helps analysts track sales activity for the fictional TICKIT web site, where users buy and sell tickets online for sporting events, shows, and concerts. In particular, analysts can identify ticket movement over time, success rates for sellers, and the best-selling events, venues, and seasons. Analysts can use this information to provide incentives to buyers and sellers who frequent the site, to attract new users, and to drive advertising and promotions.



JOINING (1)

COUNT NUMBER OF EVENTS PER CATEGORY GROUP

```
select
    catgroup,
    count(distinct eventid) as no_events
from
    sql_masterclass.event e
inner join
    sql_masterclass.category c on e.catid=c.catid
group by 1;
```

Result

select catgroup, count(distinct ever		Enter a SQL expression t
	catgroup	no_events
1	Shows	3,800
2	Concerts	4,998

```
select
    catgroup,
    count(distinct eventid) as no_events
from
    sql_masterclass.event e
right join
    sql_masterclass.category c on e.catid=c.catid
group by 1;
```

Result

	catgroup	no_events
1	Shows	3,800
2	Concerts	4,998
3	Sports	0

JOINING (2)

GET USERNAME OF THE 3 SELLERS THAT HAVE MADE MORE COMMISSION WITH THEIR SALES IN OCTOBER 2008

```
select
    username,
    sum(commission)
from
    sql_masterclass.sales s
inner join
    sql_masterclass.users u on s.sellerid=u.userid
where
    s.saletime>='2008-10-01' and s.saletime<'2008-11-01'
group by 1
order by 2 desc
limit 3;
```

Result

	username	sum
1	DVD32FRV	2,427.6
2	BEL95AGA	2,239.2
3	OSU30UYT	2,151.45

JOINING (3)

HOW MANY USERS FROM 'MEDFORD' CITY BOUGHT TICKETS FOR 'SPRING AWAKENING' EVENT

```
select
    count(distinct u.userid)
from
    sql_masterclass.event e
inner join
    sql_masterclass.sales s on e.eventid=s.eventid
inner join
    sql_masterclass.users u on u.userid=s.buyerid
where
    e.eventname='Spring Awakening' and
    u.city='Medford';
```

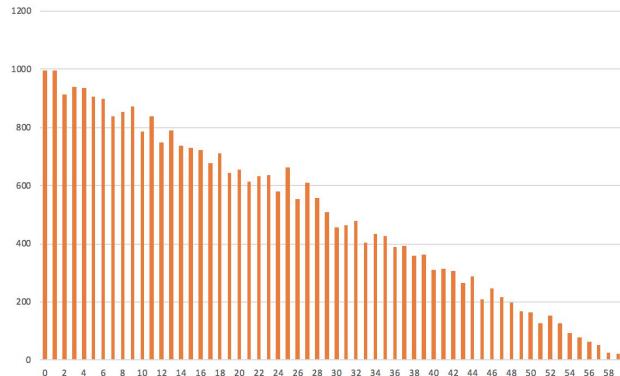
Result

select count(distinct u.userid) from sql_ma | Enter a SQL express

	count
1	5

JOINING (4)

GET DISTRIBUTION OF HOW MANY DAYS IN ADVANCE PEOPLE BOUGHT TICKETS FOR EVENTS IN VENUES IN NY STATE



```
select
    datediff(day, saletime, starttime),
    count(distinct salesid)
from
    sql_masterclass.sales s
inner join
    sql_masterclass.event e on s.eventid=e.eventid
inner join
    sql_masterclass.venue v on v.venueid=e.venueid
where
    v.venuestate='NY' and
    e.starttime='2008-07-01' and
    e.starttime<'2009-01-01'
group by 1
```

Result

date_diff	count
1	996
2	994
3	914
4	940
5	936
6	905
7	897
8	840

HOMEWORK 3

TASK:

- GET TOTAL NUMBER OF UNSOLD TICKETS PER VENUE IN OHIO STATE
- SAVE THE QUERY AS HOMEWORK3.SQL IN GITHUB IN SQL_MASTERCLASS FOLDER

TIP:

- USE SUBQUERIES TO AVOID DUPLICATION OF ROWS

```
select
    eventname
from
    sql_masterclass.event
where
    venueid IN (select distinct venueid from sql_masterclass.venue where venuestate='OH');

Result
select eventname from sql_masterclass.evi| Enter a SQL expression to filter results (use Ctrl+Space)
eventname
1 Smashing Pumpkins
2 Gogol Bordello
3 Carrie Underwood
4 Gavin DeGraw
5 Billy Joel
```





I don't
understand
at all

I need to
go over
this again

I think I got
it, but am
not
completely
comfortable

I got it

I can
explain it
to someone
else

1. INTRODUCTION TO DATABASES
2. OUR DATA WORLD
3. AMAZON REDSHIFT VS GOOGLE BIGQUERY
4. HOW TO CREATE A DATABASE
5. BASIC SQL
6. JOIN-ING TABLES
7. WINDOW FUNCTIONS
8. HOW TO EXPORT AND IMPORT DATA IN REDSHIFT
9. QA

SQL MasterClass
Barcelona Data Institute



WINDOW FUNCTIONS IN REDSHIFT

AGGREGATE FUNCTIONS

- AVG
- COUNT
- CUME_DIST
- FIRST_VALUE
- LAG
- LAST_VALUE
- LEAD
- MAX
- MEDIAN
- MIN
- NTH_VALUE
- PERCENTILE_CONT
- PERCENTILE_DISC
- RATIO_TO_REPORT
- STDDEV_POP
- STDDEV_SAMP (synonym for STDDEV)
- SUM
- VAR_POP
- VAR_SAMP (synonym for VARIANCE)

RANKING FUNCTIONS

- DENSE_RANK
- NTILE
- PERCENT_RANK
- RANK
- ROW_NUMBER

WOW, I DIDN'T KNOW SQL COULD DO THAT

WINDOW FUNCTIONS ARE A SPECIAL CLASS OF ANALYTIC FUNCTIONS THAT ARE APPLIED TO WINDOWS OF ROWS

function (expression) OVER (

[PARTITION BY expr_list]

[ORDER BY order_list [frame_clause]])

ADDING A NEW TABLES TO OUR DATASET

```
create table sql_masterclass.winsales(
salesid int,
dateid date,
sellerid int,
buyerid char(10),
qty int,
qty_shipped int);

insert into sql_masterclass.winsales values
(30001, '8/2/2003', 3, 'b', 10, 10),
(10001, '12/24/2003', 1, 'c', 10, 10),
(10005, '12/24/2003', 1, 'a', 30, null),
(40001, '1/9/2004', 4, 'a', 40, null),
(10006, '1/18/2004', 1, 'c', 10, null),
(20001, '2/12/2004', 2, 'b', 20, 20),
(40005, '2/12/2004', 4, 'a', 10, 10),
(20002, '2/16/2004', 2, 'c', 20, 20),
(30003, '4/18/2004', 3, 'b', 15, null),
(30004, '4/18/2004', 3, 'b', 20, null),
(30007, '9/7/2004', 3, 'c', 30, null);
```

SALESID	DATEID	SELLERID	BUYERID	QTY	QTY_SHIPPED
30001	8/2/2003	3	B	10	10
10001	12/24/2003	1	C	10	10
10005	12/24/2003	1	A	30	
40001	1/9/2004	4	A	40	
10006	1/18/2004	1	C	10	
20001	2/12/2004	2	B	20	20
40005	2/12/2004	4	A	10	10
20002	2/16/2004	2	C	20	20
30003	4/18/2004	3	B	15	
30004	4/18/2004	3	B	20	
30007	9/7/2004	3	C	30	

SUM

ADD AN EXTRA COLUMN TO "WINSALES" TABLE WITH TOTAL "QTY" PER SELLER

```
select
    w.*,
    s.qty_by_seller
from
    sal_masterclass.winsales w
inner join
    (select
        sellerid,
        sum(qty) qty_by_seller
    from
        sal_masterclass.winsales
    group by 1) s on w.sellerid=s.sellerid;
```



Result							
select *, sum(qty) over (partition by sellerid) Enter a SQL expression to filter results (use Ctrl+Space)							
	salesid	dateid	sellerid	buyerid	qty	qty_shipped	qty_by_seller
1	10,006	2004-01-18	1 c		10	[NULL]	50
2	10,005	2003-12-24	1 a		30	[NULL]	50
3	10,001	2003-12-24	1 c		10	10	50
4	20,002	2004-02-16	2 c		20	20	40
5	20,001	2004-02-12	2 b		20	20	40
6	30,001	2003-08-02	3 b		10	10	75
7	30,003	2004-04-18	3 b		15	[NULL]	75
8	30,007	2004-09-07	3 c		30	[NULL]	75
9	30,004	2004-04-18	3 b		20	[NULL]	75
10	40,001	2004-01-09	4 a		40	[NULL]	50
11	40,005	2004-02-12	4 a		10	10	50

```
select
    *,
    sum(qty) over (partition by sellerid) as qty_by_seller
from
    sql_masterclass.winsales;
```



	salesid	dateid	sellerid	buyerid	qty	qty_shipped	qty_by_seller
1	10,006	2004-01-18	1 c		10	[NULL]	50
2	10,005	2003-12-24	1 a		30	[NULL]	50
3	10,001	2003-12-24	1 c		10	10	50
4	20,002	2004-02-16	2 c		20	20	40
5	20,001	2004-02-12	2 b		20	20	40
6	30,001	2003-08-02	3 b		10	10	75
7	30,003	2004-04-18	3 b		15	[NULL]	75
8	30,007	2004-09-07	3 c		30	[NULL]	75
9	30,004	2004-04-18	3 b		20	[NULL]	75
10	40,001	2004-01-09	4 a		40	[NULL]	50
11	40,005	2004-02-12	4 a		10	10	50

FIRST VALUE & LAST VALUE

ADD AN EXTRA COLUMN TO "WINSALES" WITH THE FIRST "SALEID" PER SELLER

```
select
  *,
  first_value(saleid) over (partition by sellerid order by dateid)
from
  sql_masterclass.winsales;
```

FRAME CLAUSE DOES NOT APPLY TO RANKING FUNCTIONS AND IS
NOT REQUIRED WHEN NO ORDER BY CLAUSE IS USED IN THE OVER
CLAUSE FOR AN AGGREGATE FUNCTION



Query execution failed

Reason:

[Amazon](500310) Invalid operation: Aggregate window functions with an ORDER BY clause require a frame clause;

IF AN ORDER BY CLAUSE IS USED FOR AN
AGGREGATE FUNCTION, AN EXPLICIT
FRAME CLAUSE IS REQUIRED

FRAME CLAUSE IN WINDOW FUNCTIONS

IT PROVIDES THE ABILITY TO INCLUDE OR EXCLUDE SETS OF ROWS WITHIN THE ORDERED RESULT

{ ROWS | RANGE } BETWEEN *window_start* AND *window_end*

window_start

= { UNBOUNDED PRECEDING
| CURRENT ROW
| *constant-value* { PRECEDING | FOLLOWING } }

window_end

= { UNBOUNDED FOLLOWING
| CURRENT ROW
| *constant-value* { PRECEDING | FOLLOWING } }

FIRST VALUE & LAST VALUE (II)

ADD AN EXTRA COLUMN TO "WINSALES" WITH THE FIRST "SALEID" PER SELLER

```
select
  *,
  first_value(salesid) over (partition by sellerid order by dateid rows between unbounded preceding and unbounded following) as first_sale_per_seller
from
  sql_masterclass.winsales;
```

Result

select *, first_value(salesid) over (partition | Enter a SQL expression to filter results (use Ctrl+Space)

	salesid	dateid	sellerid	buyerid	qty	qty_shipped	first_sale_per_seller
1	20,001	2004-02-12	2	b	20	20	20,001
2	20,002	2004-02-16	2	c	20	20	20,001
3	40,001	2004-01-09	4	a	40	[NULL]	40,001
4	40,005	2004-02-12	4	a	10	10	40,001
5	10,001	2003-12-24	1	c	10	10	10,001
6	10,005	2003-12-24	1	a	30	[NULL]	10,001
7	10,006	2004-01-18	1	c	10	[NULL]	10,001

FIRST VALUE & LAST VALUE (III)

USING TRANSACTIONS FROM "TANGARANA" SCHEMA, GET AMOUNT IN EUROS OF THE LAST PAYMENTS MADE PER COUNTRY

```
select
    distinct user_acquisition_country,
    last_value(transaction_amount*xrate_eur)
    over (partition by user_acquisition_country order by transaction_timestamp
          rows between unbounded preceding and unbounded following) as last_amount_eur
from
    tangarana.payment_transactions t
inner join
    tangarana.users u on t.user_id=u.user_id
where
    t.transaction_type='payment';
```

Result

	user_acquisition_country	last_amount_eur
1	BM	41.761895
2	CI	39.336384
3	HT	42.31528
4	JO	42.01652
5	MU	27.863442
6	RU	33.044922
7	UM	48.199675
8	IE	16.979806
9	IR	40.70976
10	KY	27.863442

NO TRANSPARENCY NO CONSENSUS



LET'S DOUBLE CHECK...

PEER REVIEWS ARE KEY

```
select
    transaction_timestamp,
    transaction_amount,
    xrate_eur
from
    tangarana.payment_transactions
where
    transaction_type='payment' and
    user_id IN (select user_id from tangarana.users where user_acquisition_country='BM')
order by 1 desc;
```

Result

	transaction_timestamp	transaction_amount	xrate_eur
1	2017-11-29 21:50:52	33.95	1.2301
2	2017-11-21 21:28:12	33.95	1.2409
3	2017-06-12 04:47:46	33.95	1.206
4	2017-05-27 17:54:59	33.95	1.1808
5	2017-05-19 14:12:22	80.95	1.1697

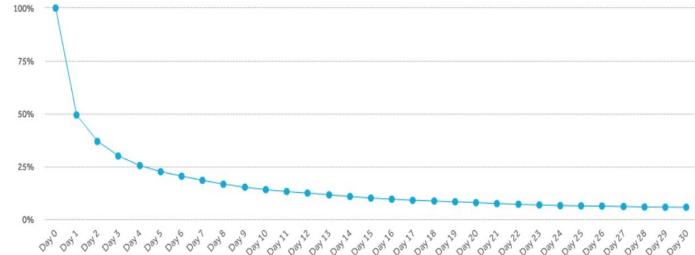
HOMEWORK 4

TASK:

- CALCULATE 30 DAYS DAILY RETENTION FOR THOSE WHO TRANSACT FOR THE FIRST TIME IN THE FIRST HALF OF 2017
- TRANSACTION RETENTION MEASURES HOW OFTEN USERS COME BACK TO TRANSACT AFTER THE FIRST TRANSACTION
- SAVE THE QUERY AS HOMEWORK4.SQL IN GITHUB IN SQL_MASTERCLASS FOLDER

TIPS:

- GET FIRST USERS THAT TRANSACT FOR THE FIRST TIME IN THE FIRST HALF OF 2017
- OUT OF THEM, CALCULATE THE VOLUME OF THOSE USERS THAT TRANSACT ON DAY0, DAY2, DAY3, ... AS SHOWN ON SCREENSHOT ON THE RIGHT
- USE THE FUNCTION DATEDIFF TO CALCULATE THE DIFFERENCE BETWEEN DATES



	datediff	count
1	0	14,925
2	1	2,198
3	2	1,253
4	3	1,315
5	4	1,295
6	5	785
7	6	553
8	7	493
9	8	377
10	9	381



I don't
understand
at all

I need to
go over
this again

I think I got
it, but am
not
completely
comfortable

I got it

I can
explain it
to someone
else

1. INTRODUCTION TO DATABASES
2. OUR DATA WORLD
3. AMAZON REDSHIFT VS GOOGLE BIGQUERY
4. HOW TO CREATE A DATABASE
5. BASIC SQL
6. JOIN-ING TABLES
7. WINDOW FUNCTIONS
8. HOW TO EXPORT AND IMPORT DATA IN REDSHIFT
9. QA

SQL MasterClass
Barcelona Data Institute

LOADING DATA IN REDSHIFT

<https://catalog.data.gov/dataset/2010-census-populations-by-zip-code>

1. GET THE DATA YOU WANT TO LOAD
2. UPLOAD DATA FILE INTO AWS S3
3. CREATE TABLE IN REDSHIFT TO STORE THAT DATA
4. LOAD DATA USING COPY COMMAND

The screenshot shows the Data.gov homepage with the URL <https://catalog.data.gov/dataset/2010-census-populations-by-zip-code>. The page title is "DATA CATALOG". The dataset title is "2010 Census Populations by Zip Code". The dataset is categorized under "City of Los Angeles / data.lacity.org". A note states: "This is a Non-Federal dataset covered by different Terms of Use than Data.gov." The dataset was last updated on February 3, 2018. The description notes: "This data comes from the 2010 Census Profile of General Population and Housing Characteristics. Zip codes are limited to those that fall at least partially within LA city boundaries. The dataset will be updated after the next census in 2020. To view all possible columns and access the data directly, visit http://factfinder.census.gov/faces/affhelp/jsf/pages/metadata.xhtml?lang=en&type=table&id=table.en.DEC_10_SF1_SF1DP1#main_content".

LOADING DATA IN REDSHIFT (II)

1. GET THE DATA YOU WANT TO LOAD
2. UPLOAD DATA FILE INTO AWS S3
3. CREATE TABLE IN REDSHIFT TO STORE THAT DATA
4. LOAD DATA USING COPY COMMAND



Amazon S3 > sql-masterclass

Overview	Properties	Permissions	Management
<input type="text"/> Type a prefix and press Enter to search. Press ESC to clear.	<input type="button" value="Upload"/> <input type="button" value="Create folder"/> <input type="button" value="Download"/> <input type="button" value="Actions"/>	EU (Ireland)	Viewing 1 to 1
Name	Last modified	Size	Storage class
<input type="checkbox"/> 2010_Census.csv	Nov 21, 2018 1:48:50 AM GMT+0100	12.2 KB	Standard

LOADING DATA IN REDSHIFT (III)

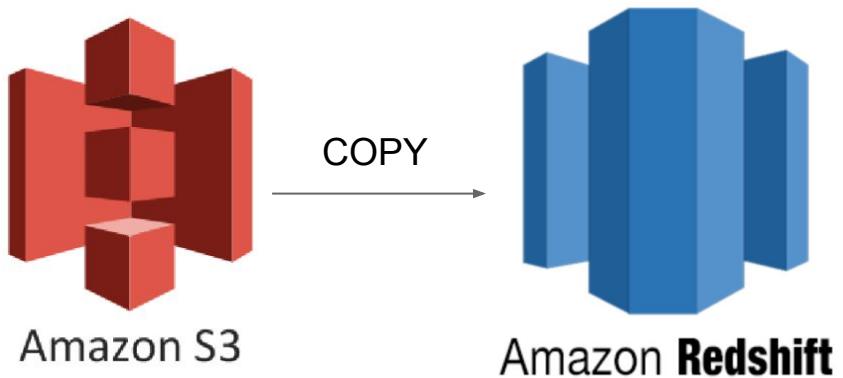
1. GET THE DATA YOU WANT TO LOAD
2. UPLOAD DATA FILE INTO AWS S3
3. CREATE TABLE IN REDSHIFT TO STORE THAT DATA
4. LOAD DATA USING COPY COMMAND



```
create table sql_masterclass.census(  
    "Zip Code" integer not null distkey sortkey,  
    "Total Population" integer,  
    "Median Age" float,  
    "Total Males" integer,  
    "Total Females" integer,  
    "Total Households" integer,  
    "Average Household Size" float);
```

LOADING DATA IN REDSHIFT (IV)

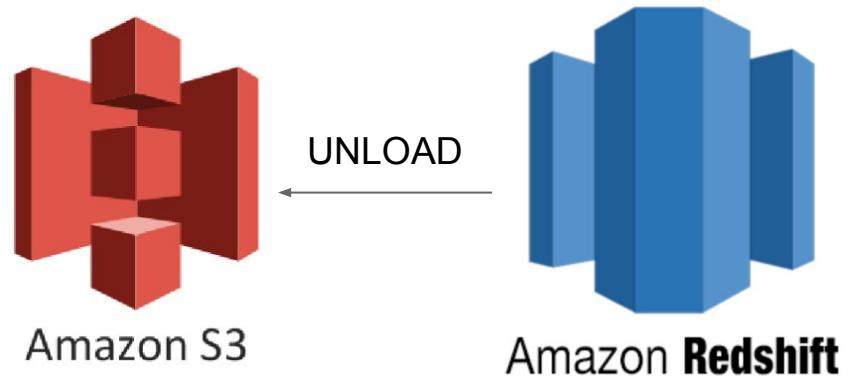
1. GET THE DATA YOU WANT TO LOAD
2. UPLOAD DATA FILE INTO AWS S3
3. CREATE TABLE IN REDSHIFT TO STORE THAT DATA
4. LOAD DATA USING COPY COMMAND



```
copy sql_masterclass.census from 's3://sql-masterclass/2010_Census.csv'  
credentials 'aws_iam_role=arn:aws:iam::529806788512:role/RedshiftCopyUnload'  
delimiter ',' IGNOREHEADER 1;
```

UNLOADING DATA FROM REDSHIFT

1. SAVE QUERY RESULTS INTO S3 BY
USING UNLOAD COMMAND



```
unload ('select
    sum("total population") as total_population,
    sum("total males") as total_males,
    sum("total females") as total_females
from
    sql_masterclass.census')
to 's3://sql-masterclass/unload_census.csv'
iam_role 'aws_iam_role=arn:aws:iam::529806788512:role/RedshiftCopyUnload'
DELIMITER as ','
PARALLEL OFF;
```

HOMEWORK 5

TASK:

- WE HAVE DATA FROM VISITS AND PURCHASES FROM AN ECOMMERCE WEBSITE IN THE FOLLOWING FORMAT:
 - "SQL_MASTERCLASS.TRANSACTIONS"
 - "SQL_MASTERCLASS.VISITS"
- ASSIGN A VISIT TO EVERY TRANSACTION
- SAVE THE QUERY AS HOMEWORK5.SQL IN GITHUB IN SQL_MASTERCLASS FOLDER

TIPS:

- USE LEAD OR LAG WINDOW FUNCTIONS TO ADD A CALCULATED COLUMN IN ONE OF THE TABLES
- LEARN HOW TO USE THE FUNCTION NVL IN REDSHIFT

The screenshot shows two separate Redshift query results. The first result is for the 'transactions' table, which has columns transaction_id, user_id, and transaction_time. The second result is for the 'visits' table, which has columns visit_id, user_id, and visit_time. Both queries return 5 rows of data.

	transaction_id	user_id	transaction_time
1	8,788,759	672,040	2016-10-03 13:52:45
2	8,708,891	649,612	2016-10-03 17:29:26
3	8,722,370	652,175	2016-10-04 04:34:10
4	8,621,989	985,347	2016-10-04 10:43:12
5	8,712,478	1,289,767	2016-10-05 09:35:57

	visit_id	user_id	visit_time
1	9,994	672,040	2016-10-03 13:28:45
2	432	1,289,767	2016-10-05 09:11:57
3	5,184	665,882	2016-10-10 20:20:17
4	8,053	1,064,435	2016-10-14 09:38:14
5	9,020	666,534	2016-10-17 11:23:00

The screenshot shows a single Redshift query result that joins the 'transactions' and 'visits' tables. The resulting table has columns user_id, transaction_id, visit_id, transaction_time, and visit_time. The query returns 10 rows of data.

	user_id	transaction_id	visit_id	transaction_time	visit_time
1	652,910	8,727,907	7,949	2016-10-06 01:33:50	2016-10-06 01:09:50
2	657,037	8,564,577	6,402	2016-10-08 12:59:28	2016-10-08 12:35:28
3	657,037	8,767,214	7,357	2016-12-24 15:01:03	2016-12-02 22:37:03
4	665,882	8,531,769	5,184	2016-10-10 20:44:17	2016-10-10 20:20:17
5	666,302	8,506,137	6,263	2016-10-05 18:48:17	2016-10-05 18:24:17
6	666,534	8,700,906	9,020	2016-10-17 11:47:00	2016-10-17 11:23:00
7	669,008	8,563,949	1,996	2016-10-17 00:33:24	2016-10-15 09:49:24
8	671,735	8,520,433	6,435	2016-12-09 05:06:32	2016-12-07 14:22:32
9	671,998	8,710,930	7,796	2016-10-23 08:26:36	2016-10-23 08:02:36
10	672,640	8,481,493	2,959	2016-10-21 03:18:29	2016-10-19 12:34:29



I don't
understand
at all

I need to
go over
this again

I think I got
it, but am
not
completely
comfortable

I got it

I can
explain it
to someone
else

1. INTRODUCTION TO DATABASES
2. OUR DATA WORLD
3. AMAZON REDSHIFT VS GOOGLE BIGQUERY
4. HOW TO CREATE A DATABASE
5. BASIC SQL
6. JOIN-ING TABLES
7. WINDOW FUNCTIONS
8. HOW TO EXPORT AND IMPORT DATA IN REDSHIFT
9. QA

SQL MasterClass
Barcelona Data Institute





**YOUR
FEEDBACK
MATTERS!**

DEEPEN LEARNING

- WITH YOUR SUPER PEN AND IN ONE POST-IT, WRITE (2 MIN):
 - YOUR 2 MAIN LEARNINGS/TAKEAWAYS
 - YOUR 2 MAIN CHALLENGES FOR THE WEEK
- FIND A PARTNER THAT IS NOT YOUR TABLE MATE
- SHARE AND HELP/SUPPORT EACH OTHER (3 MIN)
- FIND A NEW PARTNER, REPEAT (3 MIN)
 - LOOK AT WHAT IS DIFFERENT, AFTER THE FIRST SHARE



Data Analyst Progress



12%

THANK YOU

SQL MasterClass
Barcelona Data Institute