DALASOEFTY

The premiere data science training for professionals

"One should look for what is and not what he thinks should be." - Albert Einstein

INTRODUCTION TO SQL













Welcome to Intro to SQL

- Please confirm that the following technology is working properly:
 - 1. Confirm that you have SQL Server software installed correctly
 - 2. Save the following files to your desktop for this section (files previously sent to you from Data Society)
 - Claims.txt
 - crime_incidents_2013_data.csv
 - crime_incidents_2013_location.csv
 - Divisions.txt
 - Franchises.txt
 - Teams.txt

INTRODUCTION TO SQL

Overview of SQL

- 2. Working with data using SQL statements
- 3. Manipulating tables using SQL
- Logical and mathematical operations and functions 4.
- 5. SQL Server best practices

Outline





Objectives

- Develop a basic understanding of relational databases and SQL 1.
- Setting up Management Studio and importing data 2.
- 3. Learn foundational SQL query concepts





What is SQL?

- SQL is short for Structured Query Language. It is the standard language used to communicate with most relational databases
 - SQL queries are sent to a database to ask it to perform a specific task with the data it stores
- Microsoft SQL Server is a type of relational database

INTRODUCTION TO SQL







- Learn to Import/Export, manipulate, combine, and aggregate data sets using SQL Management Studio
- These lessons are structured with the intent of using SQL server as an intermediary to store and aggregate data as displayed in the image below



INTRODUCTION TO SQL

Goals for Intro to SQL



Why use SQL?

Great Average Limited

	SQL	Access	Excel
Data Size Limits	Best option for analyzing large data sets (over 1 million records)	Handles larger data sets than Excel, but can be limited by memory and space of local computers	Limited to 1,048,576 records for data se
Manipulating Data	Multiple queries can easily be combined to coerce data from multiple data sets	Querying capabilities similar to SQL with less flexibility and capabilities	Combining data sets can be difficult and prone to manual and formula errors
Analyzing / Reporting Data	Limited built-in analysis functions and lack of built-in reporting and visualization capabilities	Built in form and report capabilities for easy reporting, but more limited set of analysis functions	Many built in analysis functions, visualizations, and formatting for easy modeling and reporting
Speed	Faster Processing	Slower processing	Slower Processing
Compatibility	Compatible with most visualization, business intelligence, and statistical analysis platforms	Compatible with many visualization, business intelligence, and statistical analysis platforms	Compatible with most visualization, bus intelligence, and statistical analysis plat
Quality Control	SQL Queries create repeatable and auditable analysis processes that can be clearly commented	Access Queries also create repeatable processes, but are often less transparent than SQL queries	Excel Analyses are harder to replicate de manual steps that can be difficult to au
Learning Curve	Can be easy to learn for people without a programming background	Can be easy to learn for people without a programming background	Easy to learn for any analyst







Intro to relational databases

- What is a relational database?
 - based on common attributes in the columns and rows of those tables
 - tables



INTRODUCTION TO SQL

- Relational databases store data in the form of tables that can be related to one another - SQL Queries can leverage these relationships to rearrange the data stored in database



SQL Server general functionality

- capabilities of SQL Server:
 - Import/Store/Export data
 - Execute queries
 - Combine existing data
 - Update existing data
 - Extract, summarize, and/or aggregate data •
 - Perform calculations on existing data
 - Create and delete data tables
 - Administrative functions
 - Database security and access control

INTRODUCTION TO SQL

SQL Server Management Studio is an environment used by analysts to access the



*Screen Shot from SQL Server Management Studio

SQL Server general components

- Server database servers are programs that provides database services to other computer programs
- same fashion.
- Table data stored in a tabular format with rows of named columns
- accessing servers, databases, and tables to launch commands to the server)

INTRODUCTION TO SQL

SQL Server Management Studio				
SQLQuery1.sql - DATASOCIE	TYSERVER01.master (DS\JDOE (52))* - M			
Query Project Debug Tools Window Help				
🚺 📃 New Query 📑 📸 📸 🍒 🛓 🛋 🕰 🔊 🔸	(" - ↓ ↓ ↓ ↓			
🚽 🕴 Execute 🕨 Debug 🔲 🖌 👬	9 🔒 🎦 🧌 🥘 🦉 🏐 🗏 🗄 🛊 🛊 🐔			
- ₽ ×	SQLQuery1.sql - DATS\JDOE(102)			
🔲 🍸 🛃 🍒	1 SELECT * FROM TEST Query			
TYSERVER01 (SQL Server 12.3.6234 – DS\JDOE)				
s				
m Databases				
oase Snapshots				
BASE_01				
BASE_02				
BASE_01				
atabase Diagrams				
ables				
System Tables				
FileTables				
dbo.tbl_Data_Society_Stats				

Database – is a container of data/information organized into tables (and other structures) so that they can be easily managed and accessed back in

SQL Server Management Studio – An application used to configure, manage, and administer components of SQL server (i.e. the user interface for

Exercise 1 - connect to SQL

• In this activity we will **①** connect to a SQL Server, **②** create a new database, **B**open a new query window, and **A** select a SQL database to query from

eļ	Connect to Server ×
SQL Se	erver 2012
Server type: Server name: Authentication: Login: Password:	Database Engine DATASOCIETYSERVER01 SQL Server Authentication Remember password
Conn	ect Cancel Help Options >>

Connect to the "XXXXXXX" Server created on your computer when SQL Express was installed

Create a database and call it "Data_Society_SQL_Class"

INTRODUCTION TO SQL

Getting Started

DATA SOCIETY © 2017

• Create a new SQL script by clicking on "New Query in the toolbar

dropdown menu

10

other methods:

Import/Export Wizard Description Gives users a point and click interface Hand to upload tables into SQL server while can b using Management Studio comp Advantages Easy user interface • Cr Best use is for small tables with little • Do chance of import error: im Basic look up table containing a value and description - Sample data sets Can alter data to fit into a defined **Drawbacks** • M(table rather than resulting in an error

INTRODUCTION TO SQL

Data loading

• There are 3 methods for importing data into SQL server. This class will focus on using the Import/Export Wizard; however, it is important to understand the

SQL Server Integration Services (SSIS)
Uses Microsoft Visual Studio to crea process flow for importing and transforming tables
 Creates a repeatable audit trail Dynamically import multiple files Import similar files using loops Perform additional transformation and queries before and after table imports
 Requires learning Microsoft Visua Studio More difficult to set up and has some situations requiring new syr

Exercise 2a - loading data

- Uploading Files using the Import/Export Wizard
 - accompanying these slides
 - Claims.txt
 - crime incidents 2013 data.csv
 - crime incidents 2013 location.csv

Open Import/Export Wizard

INTRODUCTION TO SQL

- Upload 3 files using Import/Export Wizard following the directions in the SQL script

Follow Import Steps

- Use code to create and import table using bulk insert
 - accompanying these slides
 - Franchises.txt
 - Teams.txt •
 - Divisions.txt

--Step 1: Create Franchises Table _____ USE [Data Society SQL CLass] CREATE TABLE Data_Society_SQL_CLass.[dbo].[tblDemosFranchises]() ON [PRIMARY] GO --data has to be on same server copy to it first Bulk insert dbo.tblDemosFranchises From 'C:\Users\W530\Desktop\SQL Class\2a) Activity Files\Franchises.txt' With (Fieldterminator = ' ') --(120 row(s) affected)

INTRODUCTION TO SQL

Exercise 2b - loading data

- Create a table and import the following text files into SQL Server using the SQL script

[FranchID] [char](3) COLLATE SQL_Latin1_General_CP1_CI_AS NULL, [FranchName] [varchar](50) COLLATE SQL_Latin1_General_CP1_CI_AS NULL, [Active] [varchar](4) COLLATE SQL_Latin1_General_CP1_CI_AS NULL

--The console message should reflect 120 records affected if this is performed successfully

Post import data treatment

- Often tables imported into SQL are not ready for analysis
- When preparing for an analysis think of three types of tables:
 - Raw freshly imported into the SQL environment
 - Intermediate mirror of raw tables with fields transformed into usable data types Ex: dates might be imported as text and require transformation into a date format for analysis
- - Analysis tables with the fields required for an analysis

• This will be covered further in the Best Practices section

INTRODUCTION TO SQL

- *Note: it is a best practice to name tables to distinguish these table types (Ex. tbl_**Raw**_Procurement_Data _2016)
 - DATA SOCIETY © 2017

Query examples and explanation

- SQL Queries follow a standard order of statements that must
 - commands must be maintained
- 1. SELECT
- $2 \cdot INTO$
- 3. FROM
- 4. WHERE
- 5. GROUP BY
- 6. HAVING
- 7. ORDER BY

DATA SOCIETY © 2017

be followed in each query for SQL server to understand a query - Not all statements are required for every query, but the same order of

<pre>1 SELECT Contract_N0 2 FROM [tbl_DS_Sample] 3 4 SQLQuery1.sql - USS\mfergusson (53))* × Slide examp 1 □SELECT Contract_N0 2 FROM [tbl_DS_Sample] 3 WHERE Contract_N0 like '%12345F%' 4 ORDER BY Contract_N0 5 SELECT Contract_N0 like '%12345F%' GROUP BY Contract_N0 = '12345F' ORDER BY Contract_N0</pre>	LQuery	1.sql - not connected	d* 🗙 Slide		
4 SQLQuery1.sql - USS\mfergusson (53))* × Slide examp 1 □ SELECT Contract NO 2 FROM [tbl DS Sample] 3 WHERE Contract NO like '%12345F%' GROUP BY Contract NO 5 WHERE Contract NO like '%12345F%' GROUP BY Contract NO HAVING Contract NO FY Contract NO HAVING Contract NO SY Contract NO	1 2 3	SELECT Contra FROM [tbl_DS_	ct_NO Sample]		
<pre>1 SELECT Contract NO 2 FROM [tbl DS Sample] 3 WHERE Contract NO like '%12345F%' 4 ORDER BY Contract NO 5 WHERE Contract NO like '%12345F%' GROUP BY Contract NO HAVING Contract NO = '12345F' ORDER BY Contract NO</pre>	4		SQLQuery1.	ql - USS\mfergusson (53))* 🗙	Slide examp
	qI - US SELECT FROM [WHERE GROUP HAVING ORDER	S\mfergusson (53))* Contract NO tbl DS Sample] Contract NO like BY Contract NO Contract NO = ' BY Contract NO	1 E 2 3 4 5 '%12345F%	SELECT Contract NO FROM [tbl DS Sample] WHERE Contract NO like 'S ORDER BY Contract NO	%12345F%'

Query examples and explanation

INTRODUCTION TO SQL

- **SELECT** defines the fields that will be included in the new table requested by a query from other tables and from functions
- **INTO** declares that a query will create a new table in the database
- $\underline{\mathsf{FROM}}$ defines the existing tables that a query will draw data from
- **WHERE** filters the query results based on criteria from the original tables
- **GROUP BY** aggregates query results to include distinct values of the selected fields
- HAVING Similar to WHERE, except can contain aggregate functions. HAVING clauses can reference any of the items that appear in the select list
- ORDER BY sorts the query results in order by the indicated fields

Querying basics: SELECT & FROM

- The basic component of a SQL quer used to:
 - 1. Return text
 - 2. Return the results of basic math or other operations
 - 3. Returns 1 or more fields from a table in the FROM clause

SELECT STAR - Return 793 records from all columns

SELECT *
FROM tblClaims

SELECT TOP X – Returns the first **X** Records from columns selected

SELECT TOP 20 * FROM tblClaims

INTRODUCTION TO SQL

• The basic component of a SQL query is the SELECT statement which can be

ner operations

SELECT Columns – Returns 793 records of specified columns

SELECT ClaimID, MemberID FROM tblClaims

SELECT DISTINCT – Returns 11 unique records (no duplicates) across the columns selected

SELECT DISTINCT Gender

, ProviderType
FROM tblClaims

Querying basics: WHERE

- defined in the WHERE clause
 - and pattern matching

Where Specific Column Value – returns 3 records WHERE the diagnosis is equal to "V16.3"

SELECT * FROM tblClaims WHERE Diagnosis = 'V16.3'

INTRODUCTION TO SQL

• The WHERE clause is used to filter records from a table based on logical criteria

- Logical operators will be discussed in more detail later, but this includes equality, ranges,

Where value range – returns 105 records WHERE the claim payment was between \$100 and \$200

SELECT * FROM tblClaims WHERE CAST (Paid AS MONEY) > 100 AND CAST (Paid AS MONEY) < 200

Note: the Paid column was uploaded in a text formats and therefore needs to be converted to a numeric data type in order to be compared to the desired numerical range

Querying basics: GROUP BY & HAVING

- selected
 - GROUP BY statements are different from using SELECT DISTINCT because it allows aggregation functions (e.g. counts, sums, etc.)
- aggregation functions used in conjunction with the grouping

GROUP BY with count – returns 11 unique records and the counts of occurrences for each combination of records in the original data set

SELECT Gender, ProviderType , COUNT(*) AS Counts FROM tblClaims GROUP BY Gender, ProviderType

INTRODUCTION TO SQL

• The GROUP BY clause is used to return unique records (no duplicates) across columns

• The HAVING clause can only be used in conjunction with GROUP BY. It acts as an additional WHERE clause for the results of a GROUP BY statement or the results of any

> GROUP BY with Count and HAVING criteria – returns 2 of the unique records from the original query based on additional filtering from the HAVING clause

SELECT Gender, ProviderType , COUNT(*) AS Counts FROM tblClaims GROUP BY Gender, ProviderType HAVING ProviderType = '246' AND COUNT (*) >= 5

Querying basics: aggregation

- items into a single item (or multiple line items when using "GROUP BY")
- SQL aggregation functions include:
 - MIN
 - MAX
 - SUM
 - AVG
 - COUNT

tbl_cash_balances

Group	Cash_Balance	
А	20	
В	50	
В	100	
D	120	
E	30	

FROM #tbl cash balances --Aggregating by groups SELECT [GROUP]

- GROUP BY [GROUP]

INTRODUCTION TO SQL

Aggregate functions are used to summarize data by rolling up a set of data

Querying basics: sorting data

- The ORDER BY query statement is used to sort query results
 - Query results can be sorted by one or more fields
 - Using ASC or DESC after a field explicitly sorts the results in ascending or descending order, respectively
 - When there is no explicit reference the sorting defaults to ascending

tbl_cash_balances

Group	Cash_Balance	
А	20	
В	50	
В	100	
D	120	
E	30	

--Sort Ascending SELECT [GROUP], [Cash_Balance] FROM #tbl_cash_balances ORDER BY [Cash Balance]

--Sort Descending SELECT [GROUP], [Cash_Balance] FROM #tbl_cash_balances ORDER BY [GROUP] DESC

INTRODUCTION TO SQL

Note that only columns explicitly called will be sorted

Querying basics: saving results

- Saving results: If query results need to be saved for later reference or analysis there are several methods to save those results:

 - Exporting results (save results as a .csv file or copy results into a spreadsheet) - Using the INTO clause to create a permanent or temporary table
 - Permanent Table (use "INTO TableName") Creates a permanent table that will appear in a database for all users

 - instance of SQL Server

Create Permanent Table

SELECT * INTO tblClaims copy FROM tblClaims

SELECT * INTO #tblCl FROM tblCla

INTRODUCTION TO SQL

Create Loc

• Local Temporary Table (use "INTO #TableName") – Visible only to their creators during the same connection to an instance of SQL Server as when the tables were first created or referenced • Global Temporary Table (use "INTO ##TableName") – Visible to any user and any connection after they are created, and are deleted when all users that are referencing the table disconnect from the

al Temporary Table	Create Global Temporary Table
aims_tmp_local ims	SELECT * INTO ##tblClaims_tmp_globa FROM tblClaims

Querying basics: deleting tables

- the DROP TABLE command is used to remove tables from a database
- a reversible action
 - tables)
 - It is often helpful to have a statement to drop temporary tables (and sometimes streamline updates to an analysis

Dropping Tables

DROP	TABLE	tblClaims_copy
DROP	TABLE	<pre>#tblClaims_tmp_local</pre>
DROP	TABLE	##tblClaims tmp global

INTRODUCTION TO SQL

• This command will remove a table from a database permanently and this is not

- BE EXTREMELY CAREFUL when dropping tables in a database (especially permanent

permanent analysis tables) prior to the statement creating those tables in a script to

Analysis Use Example

DROP TABLE #tblDemosFranchises SELECT * INTO #tblDemosFranchises FROM [dbo].[tblDemosFranchises] --WHERE Active = 'NA' --(50 row(s) affected) (WHERE Active = 'NA') WHERE Active = 'Y' --(60 row(s) affected) (WHERE Active = 'Y')

Querying basics: commenting code

Comments are non-executing text statements that should be used to explain queries

- In line comments – all text after "--" will not be read by SQL Server when a query is run

SELECT Contract NO -- (This Text will not affect the query)

- Block comments - all text between "/*" and "*/" will not be read by SQL Server when a query

Querying basics: formatting code

Formatting and commenting SQL code should be done in a consistent and repeatable manner to make your code easier to proofread and change

Poor formatting – lines run off the screen making code hard to read

• Good formatting – comments and code are in logical order and easily read

--Step 1) Look at all Contract numbers in scope SELECT Contract NO FROM [tbl DS Sample] WHERE Contract NO = '54345F'

- SELECT Contract NO FROM [tbl DS Sample] WHERE Contract NO = '54345F'

Exercise 3 - writing basic SQL

- Practice writing basic SQL statements
- Remember to use proper formatting and appropriate comments

INTRODUCTION TO SQL

Questions?

1. Overview of SQL

Working with data using SQL statements

- 3. Manipulating tables using SQL
- Logical and mathematical operations and functions 4.
- 5. SQL Server best practices

Outline

Objectives

1. Understand, manipulate, and alter SQL Data types

2. Learn how to change data in tables

INTRODUCTION TO SQL

SQL data types

There are many different data types in SQL Server; however, there are 3 main data type categories:

- Numeric: contains numbers and can be used in mathematical operations
- Character: contains strings of text and can be searched for words and phrases or concatenated
- Date: contains dates and/or times that are stored as number allowing date type fields to also be used in mathematical operations

Data type examples

Numeric	Character	Date
1	A123F	1/01/2000
-2,000	Coffee is a great way to start off your day	2005-07-01 00:00:00:000
\$250.35	Automobile	June 16 2013
0.0023464	Desk	Monday, January 2002

Examples of data types

Numeric Data Types

Character Data Types

- Char
- Varchar
- Nvarchar

Date Data Types

- Datetime
- Date
- Time

- Int
- Money
- Bit
- Decimal

NULLS

- "zero-length" string values (i.e. "")
 - NULL values are excluded from aggregate functions:
 - it returns a count of 2
 - relationship for combining tables

BLANK Values

NULL Values

INTRODUCTION TO SQL

• NULL values are non-existing records in a field. They are different from "blank" or

Example: when SQL counts the number of records in the ID Number column ("COUNT (ID Number)")

- NULL values do not link to one another when they are in a field being used as the

• To locate NULL values, use "IS NULL" (or "IS NOT NULL") in a WHERE clause - Example: "WHERE ID Number IS NULL" would return only row 3 in the table below

	ID_Number	Category	Purchase_Date
		BMW	2002-01-01 00:00:00.000
2	1112	Mercedea	►NULL
	NULL	NULL	1970-12-01 00:00:00.000

CAST and CONVERT functions

- type into another
- Function syntax: - CAST ([FIELD NAME] AS [DATA TYPE]) - CONVERT ([DATA TYPE], [FIELD NAME])

• The CAST and CONVERT functions explicitly convert expressions of one data

--Convert Text Date to a Datetime format SELECT CONVERT (date, '1/1/2000')

--Convert Text number to a numeric format SELECT CAST ('1.0023567' AS MONEY)

Case statements

- based on logical statements
- Types of logical operators
 - Simple (equality check)
 - Searched (expressions with additional logic such as >, <, AND, OR, etc.)

```
--Case Statement Syntax
CASE WHEN ... THEN
    [WHEN ... THEN]
    [ELSE]
END
--Case Statement Example
CASE WHEN Amount 1 > Amount 2 THEN '> Amount 1'
    WHEN Amount 1 <= Amount 2 THEN '<= Amount 1'
    ELSE 'N/A'
END
```

INTRODUCTION TO SQL

• Evaluates a list of conditions and returns one of multiple result expressions

Case statement example

#tblDemosTeams_case_ex

Yr	TeamID	Points
2001	ABC	1024
2015	ABC	910
2017	BCD	500

Yr	TeamID	Offensive Rating
2001	ABC	Great Offense
2015	ABC	Unknown
2017	BCD	Poor Offense

Results

INTRODUCTION TO SQL

- WHEN Points > 910 THEN 'Great Offense'
- WHEN Points < 910 THEN 'Poor Offense'

Note that the logic in this case statement does not address a score of exactly 910 which is why one result record is 'Unknown'

Exercise 4 - data types

- Using the tables you imported earlier, practice the syntax for:
 - CASE statements
 - CAST and CONVERT functions

INTRODUCTION TO SQL

- and change data values
- These statements include:
 - INSERT: add records to a table
 - DELETE: delete records from a table
 - TRUNCATE: delete all records in a table
 - UPDATE: change values



• Data tables can also be changed using specific statements that add, remove,



INSERT statements add rows from an input source into a table

- Different syntax structures
 - Identical tables
 - Different structures
 - Values

tbl_cash_balances

Group	Cash_Balance
А	20
В	50
В	100
D	120
E	30



--Insert into the same columns **INSERT** INTO tbl_cash_balances VALUES ('G', 90)

INSERT INTO tbl_cash_balances SELECT 'E', 90

INSERT INTO tbl_cash_balances SELECT [GROUP], [Cash Balance] FROM tbl cash balances insert

--Insert different columns **INSERT** INTO tbl_cash_balances ([GROUP]) SELECT [GROUP] --with a NULL value

INTRODUCTION TO SQL



Changing data: INSERT statement





Changing data: DELETE Statement

DELETE statements permanently remove rows from data tables - Removes either an entire table or specified records

tbl cash balances

Group	Cash_Balance	
А	20	
В	50	
В	100	
D	120	
E	30	



--Remove all rows **DELETE** tbl_cash_balances

DELETE tbl_cash_balances WHERE [GROUP] = `E'

INTRODUCTION TO SQL





DATA SOCIETY © 2017





120

 \square



Changing data: TRUNCATE statement

TRUNCATE statements remove all rows from a table - The deletions are not logged making this statement faster than the DELETE statement when working with large data sets

tbl cash balances

Group	Cash_Balance
А	20
В	50
В	100
D	120
E	30



tbl cash balances --Remove all rows TRUNCATE TABLE tbl_cash_balances Group Cash Balance Note that TRUNCATE will leave a table with columns but no records similar to DELETE without criteria

INTRODUCTION TO SQL







Changing data: UPDATE statement

conditional logic

tbl_cash_balances

Group	Cash_Balance	
А	20	
В	50	
В	100	
D	120	
E	30	

--Remove all rows **UPDATE** tbl_cash_balances SET [Cash Balance] = 150

--Remove specific Values based on --conditional logic **UPDATE** tbl_cash_balances SET [Cash Balance] = 150WHERE [GROUP] = 'E'

INTRODUCTION TO SQL

• Updates the values either in an entire column or based on specific values using





Exercise 5 - changing data

- Using the tables you imported earlier, practice the syntax for:
 - INSERT
 - DELETE
 - TRUNCATE
 - UPDATE







INTRODUCTION TO SQL



Questions?



1. Overview of SQL

- 2. Working with data using SQL statements
- 3. Manipulating tables using SQL
- Logical and mathematical operations and functions 4.
- 5. SQL Server best practices

Outline







- Understand table combination methods
- Recognize types of joins and unions
- Understand table relationships and keys
- Creating views
- Practice writing JOIN & UNION statements in SQL

Objectives



SQL table combinations

SQL tables can be combined using JOIN or UNION statements to merge columns or records respectively

JOIN

ID	First Name	ID	Last Name	ID	First Name	Last
A1	John	A1	Smith	A1	John	S
A2	Samantha	A2	Ripken	A2	Samantha	Ri
A3	Paul	A3	Johnson	A3	Paul	Joł

- A JOIN brings columns from 2 different tables into a combined table
- The combined table can have more or fewer records than its parent tables depending on both the relationship between records on the parent tables and the type of join performed



UNION

ID	First Name	Last Name
A1	John	Smith
A2	Samantha	Ripken
		I
_		
D	First Name	Last Name
A3	Paul	Johnson
A4	Taylor	Prince

- A UNION appends records from 2 tables into a combined table
- The combined table will have more records than its parent tables assuming no filtering is applied to either table

Name mith ipken hnson





Using table aliases

A table alias allows a table to be temporarily renamed within the scope of an individual query. This makes queries easier to read and limits the amount of code required for a query to execute.





Joins: code structure

The illustration below demonstrates how the JOIN code relates these tables. In more complex joins, tables are referred to as LEFT and RIGHT tables based on their order in the SQL statement, which will impact the records returned in result sets.



(LEFT TABLE)

INTRODUCTION TO SQL





Joins: INNER JOIN

An INNER JOIN between 2 tables returns the intersection between those 2 tables



ID	First Name	Last Name
A1	John	Smith
A3	Paul	Johnson

INTRODUCTION TO SQL





Joins: LEFT OUTER JOIN

A LEFT OUTER JOIN between 2 tables returns all records from the table in the initial table and the intersection between those 2 tables. Where there is no intersection NULL values are populated in columns selected from the joined table.



Result

ID	First Name	Last Name
A1	John	Smith
A2	Samantha	NULL
A3	Paul	Johnson

FROM [First] A LEFT JOIN [Last] B ON A.ID = B.ID

INTRODUCTION TO SQL









Joins: LEFT OUTER JOIN (exclude)

A LEFT OUTER JOIN with exclusion between 2 tables returns only records from the original table with no intersection to the initial table.

Tables -

[First]

ID	First Name
A1	John
A2	Samantha
A3	Paul

ID	Last Name
A1	Smith
A3	Johnson
A4	Adler

FROM [First] A LEFT OUTER JOIN [Last] B ON A.ID = B.IDWHERE B.ID IS NULL



Result

ID	First Name	Last Name
A2	Samantha	NULL

FROM [First] A LEFT JOIN [Last] B ON A.ID = B.IDWHERE B.ID IS NULL

INTRODUCTION TO SQL







Joins: RIGHT OUTER JOIN

A RIGHT OUTER JOIN between 2 tables returns all records from the joined table and the intersection between those 2 tables. Where there is no intersection NULL values are populated in columns selected from the table in the initial table.

Tables

[First]		
ID	First Name	
A1	John	
A2	Samantha	
A3	Paul	

ID	Last Name	
A1	Smith	
A3	Johnson	
A4	Adler	

[last]

, B.[Last Name] FROM [First] A ON A.ID = B.ID



Result

ID	First Name	Last Name
A1	John	Smith
A3	Paul	Johnson
A4	NULL	Adler

, B.[Last Name]] FROM [First] A RIGHT JOIN [Last] B ON A.ID = B.ID

INTRODUCTION TO SQL







Joins: RIGHT OUTER JOIN (exclude)

A RIGHT OUTER JOIN with exclusion between 2 tables returns only records from the joined table with no intersection to the initial table.



First Name Last Name

Adler

NULL

SELECT B.ID, , B. [Last Nam FROM [First] RIGHT OUTER ON A.ID =WHERE A.ID IS

SELECT B.ID, , B.[Last Nam FROM [First] RIGHT JOIN ON A.ID = WHERE A.ID I

INTRODUCTION TO SQL

Α4



Code	– Logic —
A.[First Name] ne] A JOIN [Last] B = B.ID 5 NULL	A B
Or	
A.[First Name] me] A Last] B = B.ID S NULL	Connection NULL ↔ B





Joins: FULL OUTER JOIN

A FULL OUTER JOIN between 2 tables returns all records from both tables including their intersection.



ID	First Name
A1	John
A2	Samantha
A3	Paul

ID	Last Name
A1	Smith
A3	Johnson
A4	Adler



Result

ID	First Name	Last Name
A1	John	Smith
A2	Samantha	NULL
A3	Paul	Johnson
A4	NULL	Adler

Code Logic CASE WHEN A.ID IS NULL SELECT THEN B.ID ELSE A.ID END as ID FROM [First] A FULL OUTER JOIN [Last] B ON A.ID = B.IDConnection A↔ NULL

INTRODUCTION TO SQL



DATA SOCIETY © 2017



 $\mathsf{NULL} \rightarrow \mathsf{B}$



Joins: FULL OUTER JOIN (exclude)

A FULL OUTER JOIN with exclusion between 2 tables returns only records from both tables excluding their intersection.

Tables -

[First]

ID	First Name
A1	John
A2	Samantha
A3	Paul

ID	Last Name
A1	Smith
A3	Johnson
A4	Adler

[Last]



Result

ID	First Name	Last Name
A2	Samantha	NULL
A4	NULL	Adler

Code Logic SELECT CASE WHEN A.ID IS NULL THEN B.ID ELSE A.ID END as ID , A.[First Name] , B.[Last Name] FROM [First] A FULL OUTER JOIN [Last] B ON A.ID = B.IDConnection WHERE A.ID is NULL OR B.ID IS NULL A ↔ NULL

INTRODUCTION TO SQL

DATA SOCIETY © 2017



 $NULL \leftrightarrow B$



Joins: CROSS JOIN

to the other.



Code Logic SELECT CASE WHEN A.ID IS NULL THEN B.ID ELSE A.ID SELECT A. [First Name] , B.[Last Name] FROM [First] A CROSS JOIN [Last] B Connection

INTRODUCTION TO SQL

ID



A CROSS JOIN between 2 tables returns every combination of records from one table







Joining logic comparison







Union statements

• Unions allow records from the same fields to be appended to one another



Date	Revenue
12/10/01	\$25,000
12/19/01	\$120,000
12/23/01	\$10,000
01/10/01	\$35,000
01/19/01	\$12,000
01/23/01	\$110,000

INTRODUCTION TO SQL

Code

SELECT Date, Revenue FROM Dec_01_Sales UNION ALL SELECT Date, Revenue FROM Jan_02_Sales

Oľ

SELECT Date, Revenue FROM Dec_01_Sales UNION SELECT Date, Revenue FROM Jan_02_Sales



UNION VS UNION ALL

Using only UNION to combine tables will remove duplicate records between and within tables, whereas using UNION ALL will combine all records with no further alterations. Using UNION ALL typically runs significantly faster than a UNION unless the dataset is small and removing duplicates is necessary.

[Dec_01_Sales]

Date	Revenue
12/10/01	\$25,000
12/19/01	\$120,000
12/19/01	\$120,000
12/19/01	\$120,000
12/23/01	\$10,000

[Jan_02_Sales]

Date	Revenue
01/10/02	\$35,000
12/23/01	\$10,000
01/19/02	\$12,000
01/23/02	\$110,000



Result

Date	Revenue
12/10/01	\$25,000
12/19/01	\$120,000
12/23/01	\$10,000
01/10/01	\$35,000
01/19/01	\$12,000
01/23/01	\$110,000

INTRODUCTION TO SQL

UNION	ALL	
Dec_01_Sales]	[Jan_02	Sales

Date	Revenue
12/10/01	\$25,000
12/19/01	\$120,000
12/19/01	\$120,000
12/19/01	\$120,000
12/23/01	\$10,000

Date	Revenue
01/10/02	\$35,000
12/23/01	\$10,000
01/19/02	\$12,000
01/23/02	\$110,000





Table relationships

The previous examples use tables with one-to-one relationships; however, additional considerations should be taken when other table relationships exist. These relationships include:

- One-to-One each record in one table will have no more than one matching record in a second table, and vice versa.
- One-to-Many each record in one table can have many matching only have one matching record in the first table.
- in a second table, and vice versa.



records in a second table; however, each record in the second table can

• Many-to-Many – records in one table can have many matching records



One-to-Many relationships

numerous records in another table.



INTRODUCTION TO SQL



One-to-Many relationships exist when each records in one table may relate to

Items Purchased Account A153 Pizza A153 Bread B634 Gift Card B754 Soda B754 Sports Drink



Many-to-Many relationships

Many-to-Many relationships exist when records in one table can have many matching records in a second table, and vice versa.



INTRODUCTION TO SQL



Multiple JOIN statements

table joins is important (see ordering example on the next slide).

——Tables——							
[First]			[Last]		[PhoneNo]		
ID	First Name		ID	Last Name		ID	Phor
A1	John		A1	Smith		A1	555-111·
A2	Samantha		A3	Johnson		A2	<u>555-222</u>
A3	Paul		A4	Adler		A3	<u>555-333</u>
		I				A4	555-444

Result

ID	First Name	Last Name	Phone
A1	John	Smith	555-111-4564
A2	Samantha	NULL	555-222-6781
A3	Paul	Johnson	555-333-9754
A4	NULL	Adler	555-444-3456

• A JOIN statement is limited to two combining 2 tables, but a query can have multiple JOIN statements. Not all tables have to return a value. The order of

Code



SELECT A. [ID], ,B.[First Name] ,C.[Last Name] ,A.[Phone] FROM PhoneNo A LEFT JOIN [First] B ON A.ID = B.IDLEFT JOIN [Last] C ON A.ID = C.ID



Multiple JOIN statements: joins order

The order of joins will dictate what information is produced in the results. Understanding the relationships (or lack thereof) between data in tables becomes very important as tables are combined.



ClientNames

ID	Client
1	XYZ
2	ABD
5	MYR

InvoiceNo

ID	Invoice
1	H123
4	B346
5	MZ905

FROM ClientNames A LEFT JOIN InvoiceNo B ON A.ID = B.IDLEFT JOIN PhoneNo C

PhoneNo

Client	PhoneNo
XYZ	555-567-5585
ABD	555-675-1900
QAR	555-111-1801

FROM InvoiceNo A LEFT JOIN ClientNames B ON A.ID = B.ID LEFT JOIN PhoneNo C

INTRODUCTION TO SQL







Multi criteria joins

Joins can use more than 1 criteria as well as multiple fields to combine data from tables. The join logic can also use static values.



DS_Supply

City	State	Data Scientists
Springfield	Missouri	13
Springfield	Virginia	100
Vienna	Virginia	75
Springfield	New York	25
New York	New York	450

DS Demand

City	State	Data Scientist Demand
Springfield	Missouri	25
Springfield	Virginia	75
Vienna	Virginia	50
Springfield	New York	20
New York	New York	500

Code⁻

SELECT A.City, A.State , A. [Data Scientists] , B.[Data Scientist Demand] , 1.00 * A.[Data Scientists] / B.[Data Scientist Demand] as 'Demand Filled %' FROM DS Supply A LEFT JOIN DS Demand B ON A.City = B.City AND A.State = B.StateAND A.State = 'New York'

INTRODUCTION TO SQL





City	State	Data Scientists	Data Scientist Demand	DF
Springfield	Missouri	13	NULL	
Springfield	Virginia	100	NULL	
Vienna	Virginia	75	NULL	
Springfield	New York	25	20	
New York	New York	450	500	





Non-equivalent joins

Earlier examples illustrate equivalent joins (Value 1 = Value 2), but joins can also use non equivalent logic using operators such as >, <, >=, and <=.



ArticlePages

Article	Beginning Page	Ending Page
Golf	1	5
Philosophy	6	10
Science	11	25

Excerpts

Page	Text
4	Tiger won again
8	I think therefore I am
22	The chemical is inert

SELECT A.Article, B.Text FROM ArticlePages A LEFT JOIN Excerpts B



Code

Results

ON A. [Beginning Page] <= B.Page AND A. [Ending Page] >= B.Page

Article	Text
Golf	Tiger won again
Philosophy	I think therefore
Science	The chemical is i





Self joins

as creating cumulative sums over time.

Tables

Financials

Period	Revenue	Profit
1	\$ 10,000	\$ -500
2	\$ 12,000	\$ -1,000
3	\$ 11,500	\$ 200
4	\$ 13,000	\$ 250
5	\$ 17,500	\$ 400
6	\$ 14,000	\$ 500
7	\$ 9,000	\$ -100
8	\$ 20,500	\$ 1,000
9	\$ 12,000	\$ 1,200
10	\$ 15,000	\$ 800
11	\$ 19,000	\$ 950
12	\$ 21,000	\$ 1,250

SELECT A.Period SUM (B. Revenue) SUM(B.Profit) FROM #Financials JOIN #Financials ON A.Period >= H GROUP BY A. Period



SQL also allows a table to be joined to itself. This can be useful in situations such



Results

)	as	Cumulative_Revenue
	as	Cumulative_Profit
5	A	
5	В	
3	.Pe	riod

Period	Cumulative	Cumu
	Revenue	Pro
1	\$ 10,000	(
2	\$ 22,000	\$ -
3	\$ 33,500	\$ -
4	\$ 46,500	\$ -
5	\$ 64,000	(
6	\$ 78,000	(
7	\$ 87,000	(
8	\$ 107,500	
9	\$ 119,500	\$
10	\$ 134,500	\$
11	\$ 153,500	\$
12	\$ 174,500	\$





Quality control tips for joins

- Pay attention to the row counts!
 - Is the number of rows what you expect?
 - A common problem is assuming that a field or combination of fields represents a unique value. When that is not the case you can see an increase in record count.
 - Another common issue is incorrectly setting criteria of a join or WHERE clause and excluding more • records than was originally intended.
 - Use simple queries on each original table to compare the number of records with given criteria to the number in the joined table.
 - It is also helpful to use left and right joins Use LEFT and/or RIGHT joins to determine the overlap between tables and whether:
 - The fields you are using to link the original tables are appropriate.
 - The tables themselves are appropriate to combine for an analysis (2 tables can appear to be similar data from field names but actually contain little to no overlap).



A SQL view is a virtual table that is the stored results of an underlying SQL statement. That SQL statement can be a table or query. Views can be valuable for:

- Creating simplicity by hiding complex queries from end users of data
- Creating security through hiding fields with private information and/or preventing changes to base tables
- 3.







Preventing redundancy & increase consistency by providing a common source for data users



- Practice writing JOIN statements
- Remember to use proper formatting and appropriate comments



Exercise 6 - practice JOINs







INTRODUCTION TO SQL



Questions?



1. Overview of SQL

- 2. Working with data using SQL statements
- 3. Manipulating tables using SQL

Logical and mathematical operations and functions 4.

5. SQL Server best practices

Outline




Objectives

- Understand logical operators in SQL and how they are used
- Understand the use and syntax of various functions





Logical operators: comparisons

in CASE Statements, JOINS, WHERE clauses, HAVING clauses.

Comparison Operators:

- "=" and "<>" / "!=" "Equal to" and "Not Equal to"
- "<" and "<=" / "!>" "Less than" and "Less than or equal to"
- ">" and ">=" / "!<" "Greater than" and "Greater than or equal to"
- IN used to match values in a field to a list of values
- BETWEEN used to specify an inclusive range (lower and upper values are searched as well as the values in between)
- LIKE used to identify patterns in character fields
- NOT used to negate conditions

Logical operators test whether or not a condition is true. They are generally used



Logical operators: IN

The IN clause is used to compare a value to a list of possible values.

-Tables-

[Rainbow]

Rainbow Colors	Order
Red	1
Orange	2
Yellow	3
Green	4
Blue	5
Indigo	6
Violet	7



Favorite
Colors
Blue
Indigo
Orange
Black

Code

Example 1 - Field to Field comparison

SELECT * FROM Rainbow WHERE [Rainbow Colors] IN (SELECT [Favorite Colors] FROM Favorite)

Example 2 - Field to list comparison

SELECT * FROM Rainbow WHERE [Rainbow Colors] IN('Blue', 'Violet')

Example 3 - CASE statement using IN

SELECT *, CASE WHEN [Rainbow Colors] IN (SELECT [Favorite Colors] FROM Favorite) THEN 'Y' ELSE 'N' END as Fav Flag FROM Rainbow

INTRODUCTION TO SQL

Results

DATA SOCIETY © 2017

Example 1

Rainbow Colors	Order
Orange	2
Blue	5
Indigo	6

Example 2

Rainbow Colors	Order
Blue	5
Violet	7

Example 3

Rainbow Colors	Order	Fav_Flag
Red	1	Ν
Orange	2	Y
Yellow	3	Ν
Green	4	Ν
Blue	5	Y
Indigo	6	Y
Violet	7	Ν



Logical operators: BETWEEN

The BETWEEN clause can be used to determine if a value falls between an upper and lower bound. For SQL, the BETWEEN statement is inclusive of the upper and lower bounds (i.e. Between is true when [lower bound] <= value <= [upper bound]).

-Tables-

Code-

Example 1 - Text

SELECT * FROM Sales WHERE Category BETWEEN 'B' AND 'D'

Example 2 - Numeric

SELECT * FROM Sales WHERE Amount BETWEEN 25 AND 75

Example 3 - Dates

SELECT * FROM Sales WHERE Date BETWEEN `2015-12-25' AND `2016-01-20'

[Sales]

Date	Amount	Category
12/1/15	150	С
12/25/15	20	А
1/15/15	35	D
1/31/16	75	В



Results

Example 1

Date	Amount	Category
12/1/2015	150	С
1/15/2016	35	D
1/31/2016	75	В

Example 2

Date	Amount	Category
1/15/2016	35	D
1/31/2016	75	В

Example 3

Date	Amount	Category
12/25/2015	20	А
1/15/2016	35	D
1/31/2016	75	В



Logical operators: LIKE

understand how to leverage special characters ("%", "_", "[", "]", and "^") to accurately match patterns.

- Wildcard characters that substitute for any other character in a string • "%" - used to represent 0 or more characters and is typically used before and/or after the part of text being searched for to look for that text anywhere in the character string
- - "_" used to represent 1 character
- Specified patterns Brackets "[]" can be used to specify lists or ranges of characters or numbers that should be represented by a character in a pattern
 - Using "^" after the opening bracket changes exclude the characters following it (ex. "[^m]" matches any letter other than "m")
- Escape The special characters listed above ("%", "_", "[", and "]") need to be treated differently than others. They either need to be included in brackets or placed after an escape character



The LIKE clause matches patterns of text to a character field. It is important to



Logical operators: LIKE examples 1

-Code

Example 1 - Using '%' wildca
SELECT CASE WHEN 'aabbcc' LIKE '%bb%' THEN 'T' ELSE SELECT CASE WHEN 'aabbcc' LIKE '%a%c%' THEN 'T' ELSE
Example 2 - Using '_' wildca
SELECT CASE WHEN 'fair' LIKE '_air' THEN 'T' ELSE SELECT CASE WHEN 'lair' LIKE '_air' THEN 'T' ELSE
Example 3 - Using '[]' searchir
<pre>SELECT CASE WHEN 'theatre 7' LIKE '%theat[re][re] SELECT CASE WHEN 'theater 2' LIKE '%theat[re][re]</pre>
Example 4 - Using '[^]' searchi
<pre>SELECT CASE WHEN 'theatre X' LIKE '%theat[re][re] SELECT CASE WHEN 'theater 9' LIKE '%theat[re][re]</pre>
Example 5 – Using escapes for special o
SELECT CASE WHEN $a[-\%]$ LIKE $a[[]-[\%][_]$ THEN 'SELECT CASE WHEN $a[-\%]$ LIKE $a \in [-\%]$ 'ESCAPE '

INTRODUCTION TO SQL





Logical operators: LIKE examples 2

-Tables-

[txtLog]

ID	Note
1	He doesn't want to give
	a 50% cut
2	Cut that out we need
	the code
3	Here is the code
	[code:123_tf]
4	How long the code
	good for?
5	123-tf won't be good
	for too long
6	I think 124_tf will work
	too

Example 1 - word anywhere in string

SELECT * FROM txtLog WHERE Note LIKE '%code%'

Example 2 - word at end of string

SELECT * FROM txtLog WHERE Note LIKE 'Stoo'

Example 3 - skipping multiple characters

SELECT * FROM txtLog WHERE Note LIKE \%12%tf%'

Example 4 - skipping single characters and brackets

SELECT * FROM txtLog WHERE Note LIKE \%12 []tf%'

SELECT * FROM txtLog WHERE Note LIKE \%12 [^]tf%'

Code

Example 5 - skipping single characters and caret brackets

DATA SOCIETY © 2017

Results

Example 1

ID	Note
2	Cut that out we need the code
3	Here is the code [code:123_tf]
4	How long the code good for?

Example 2

ID	Note
6	I think 124_tf will work too

Example 3

ID	Note
3	Here is the code [code <mark>:123 tf]</mark>
5	123-tf won't be good for too long
6	I think 124 tf will work too

Example 4

ID	Note
3	Here is the code [code 123 tf]
6	l think 124 tf will work too

Example 5

ID	Note
	5 <mark>123-tf</mark> won't be good for too long





Multiple logical operators

Multiple logical operators can be strung together by relating them with AND or OR statements.

- meeting all conditions
- any of these conditions

• AND - used to connect two or more conditions and only returns those rows

• OR - used to connect two or more conditions and returns any rows that meet



Multiple logical operators examples

City	State	Data Scientists
Springfield	Missouri	13
Springfield	Virginia	100
Vienna	Virginia	75
Springfield	New York	25
New York	New York	450

OR

SELECT *	
FROM #DS_Supply	
WHERE City = 'Vienna'	
AND State = 'Virginia'	

AND

SELECT * FROM #DS Supply WHERE City = 'Vienna' OR State = 'Missouri'

City	State	Data Scientists
Vienna	Virginia	75

City	State	Data Scientists
Springfield	Missouri	13
Vienna	Virginia	75

INTRODUCTION TO SQL



DS_Supply

Combined

SELECT * FROM #DS Supply WHERE State = 'New York' AND [Data Scientists] > 50 OR State = 'Missouri'

City	State	Data Scientists
Springfield	Missouri	13
New York	New York	450



Functions

- within the context of a query.
- single column.



• Functions in SQL Server are stored programs that can be passed parameters and return a value or values. In SQL most functions are structured to work

• SQL has functions that are designed to work with each and/or multiple data types and also operate either across columns in a single row or across rows in a



Functions: cross column vs. aggregate

- Cross column functions combine data in the same record from different columns
- precise

Financials

Period	Revenue	Costs	=	Profit	
1	\$ 10,000	\$ 500	\rightarrow	\$9,500	
2	\$ 12,000	\$1,000	\rightarrow	\$11,000	
3	\$ 11,500	\$ 200		\$11,300	
4	\$ 13,000	\$ 250		\$12,750	
5	\$ 17,500	\$ 400	\rightarrow	\$17,100	
6	\$ 14,000	\$ 500	\rightarrow	\$13,500	
	SUM -				U
	Tot Revenue	Tot Costs		Tot Costs	
	\$ 78,000	\$ 2,850	C	1	

• Aggregate functions combine data in multiple rows from the same columns • The two function types can be used together; however, the syntax needs to be

> --Cross Column function: Periodic Profit SELECT Revenue - Costs AS 'Profit' FROM Financials --Aggregate function SELECT SUM (Revenue) AS 'Tot Revenue' , SUM(Costs) AS ' Tot Costs' FROM Financials --Cross Column: count if Profit > 15k SELECT SUM(CASE WHEN Revenue - Costs > 15000 THEN 1 ELSE 0 END) AS 'Periods GT 15K' FROM Financials



Functions: math

- Cross column functions:
 - Numeric data types support all basic math functions such as addition ("+"), subtraction ("-"), multiplication ("*"), division ("/"), and exponents ("POWER()")
 - Cross-column also includes mathematical expressions such as absolute value ("ABS()"), rounding ("ROUND()") and random number generation ("RAND()")
- Aggregate Functions:
 - SQL also supports aggregate math functions such as counts ("COUNT()"), sums ("SUM()"), maximum values ("MAX()"), and averages ("AVG()")



Functions: text

- Cross-column functions:
 - - Extract part of a string: LEFT(), RIGHT(), SUBSTRING()
 - Finding the position of text and length of string: PATINDEX(), CHARINDEX(), LEN()
 - **Removing spaces** from the ends of strings: LTRIM(), RTRIM(), TRIM()
 - Combine strings: CONCAT(), "+"
- Aggregate Functions:

- Text data types support a variety of functions to manipulate text strings including:

- Text data only supports several aggregate functions including counts ("COUNT()"), maximum alphanumeric values ("MAX()"), and minimum alphanumeric values ("MIN()")



Functions: dates

- Cross-column functions:
 - Date data types support a variety of functions to manipulate dates including:
 - Get current date and/or time: CURRENT_TIMESTAMP, GETDATE(), SYSDATETIME()
 - Return date and/or time parts: DAY(), MONTH(), YEAR(), DATENAME(), DATEPART()
 - Create date and/or time from parts: DATEFROMPARTS(), TIMEFROMPARTS() •
 - Date math: DATEDIFF(), DATEADD()
 - - converted back knowing that Microsoft's base date (numerically zero) is 1/1/1900.



- Date data only supports several aggregate functions including counts ("COUNT()"), maximum alphanumeric values ("MAX()"), and minimum alphanumeric values ("MIN()") • This limitation can be avoided because dates can be converted to a numeric data type for additional calculations. After calculations (averages, variances, etc.) are performed on a numeric date they can be



Functions: NULLs

- sets that contain NULL values:
 - ISNULL() replaces NULL with the specified replacement value
 - COALESCE() returns the first non-null expression among its arguments



[tblSales]

Date	SalesAmt1	SalesAmt2
12/1/15	150	75
12/25/15	NULL	50
1/15/16	35	NULL
1/31/16	NULL	65

--Manually substitute for NULL values SELECT Date, ISNULL (SalesAmt1, 0) AS SalesAmt FROM tblSales

--Take first non-NULL values SELECT Date, COALESCE (SalesAmt1, SalesAmt2) AS SalesAmt FROM tblSales

INTRODUCTION TO SQL

• There are some useful functions that assist with dealing with incomplete data

Results

ISNULL

Date	SalesAmt
12/1/15	150
12/25/15	0
1/15/15	35
1/31/16	0

COALESCE

Date	SalesAmt
12/1/15	150
12/25/15	50
1/15/15	35
1/31/16	65



Exercise 7 - logic & functions

- Using the tables you imported earlier, practice the syntax for:
 - Logical operators
 - Cross-column and aggregate functions
 - Math
 - Text
 - Date

INTRODUCTION TO SQL







INTRODUCTION TO SQL



Questions?



1. Overview of SQL

- 2. Working with data using SQL statements
- 3. Manipulating tables using SQL
- Logical and mathematical operations and functions 4.

SQL Server best practices

Outline





Objectives

- Understand common data problems
- Review Quality Control (QC) tips
- Review Organizational tips
- Understand ways to improve performance and efficiency



common data problems

- Unreliable or unusable data
- Incomplete data CCN 13009976 – Duplicate data 13009981 13009981 Inconsistent formats 13009981 Inaccurate interpretations NULL NULL - What do the values in OFFENSE 13009990 and METHOD really mean? 13009991 13009994 13010012 13010013 Join issues 13010014
 - NULLs not accounted for
 - Wrong columns used
 - Table relationships not considered

INTRODUCTION TO SQL

REPORTEDDATE	REPORTEDTIME	SHIFT	OFFENSE	METH
1/23/2013	NULL	DAY	NULL	OTHE
1/23/2013	9:46:00 AM	DAY	MOTOR VEHICLE THEFT	OTHE
1/23/2013	9:46:00 AM	DAY	MOTOR VEHICLE THEFT	OTHE
1/23/2013	9:46:00 AM	DAY	MOTOR VEHICLE THEFT	OTHE
13009986	1/23/2013	11:25:00 AM	DAY	THEFT F/
13009986 13009987	1/23/2013 1/23/2013	11:25:00 AM 12:30:00 PM	DAY DAY	THEFT F/ THEFT/C
13009986 13009987 01232013	1/23/2013 1/23/2013 10:54:00 AM	11:25:00 AM 12:30:00 PM DAY	DAY DAY MOTOR VEHICLE THEFT	THEFT F/ THEFT/C OTHE
13009986 13009987 01232013 2013-1-23	1/23/2013 1/23/2013 10:54:00 AM 23:57:00	11:25:00 AM 12:30:00 PM DAY DAY	DAY DAY MOTOR VEHICLE THEFT BURGLARY	THEFT F/ THEFT/C OTHE OTHE
13009986 13009987 01232013 2013-1-23 2013-1-23	1/23/2013 1/23/2013 10:54:00 AM 23:57:00 12:13:00 PM	11:25:00 AM 12:30:00 PM DAY DAY DAY	DAY DAY MOTOR VEHICLE THEFT BURGLARY THEFT/OTHER	THEFT F/ THEFT/C OTHE OTHE KNIF
13009986 13009987 01232013 2013-1-23 2013-1-23 1/23/2013	1/23/2013 1/23/2013 10:54:00 AM 23:57:00 12:13:00 PM 18:34:00	11:25:00 AM 12:30:00 PM DAY DAY DAY DAY	DAY DAY MOTOR VEHICLE THEFT BURGLARY THEFT/OTHER THEFT F/AUTO	THEFT F/ THEFT/C OTHE OTHE KNIF OTHE
13009986 13009987 01232013 2013-1-23 2013-1-23 1/23/2013 1/23/2013	1/23/2013 1/23/2013 10:54:00 AM 23:57:00 12:13:00 PM 18:34:00 12:31:00 PM	11:25:00 AM 12:30:00 PM DAY DAY DAY DAY DAY	DAY DAY MOTOR VEHICLE THEFT BURGLARY THEFT/OTHER THEFT F/AUTO THEFT/OTHER	THEFT F/ THEFT/C OTHE OTHE KNIF OTHE KNIF





QC: preventing data problems

problems:

- Reconcile data source totals
 - Compare record counts and sum totals of data sets used to source documents
 - Reconcile with other source documents
 - Compare detail and summary data sets
 - Account for all records
- Perform frequency distributions on the data
 - Check for incorrectly excluded or included data
 - Check for invalid entries
 - Distinguish between NULL and blank fields



The following are quality control (QC) steps to detect and prevent common data



QC: preventing data problems

- Check for reasonableness
 - Check that fields provide the correct data for your analysis
 - implied decimal)

 - totals or the meaning of the data
- Lookup table verification
 - Determine valid values list for lookup tables
 - Compare valid values to data values used
 - Identify missing values
- Query result tests
 - Evaluate whether or not JOIN results meet expectations
 - Verify the type of join used was correct
 - Determine if results contained too many or few records

INTRODUCTION TO SQL



- Check that amounts fields are appropriate (i.e. 5000 could be \$5,000 or \$50 if there is an

- Ensure that fields are logically related to each other and carefully compare fields for JOINS - In transactional data check for adjustments/corrections or reverse values that may impact



QC: building quality into analysis

The diagram below illustrates that quality control (QC) must be built into an analysis from beginning to end.



INTRODUCTION TO SQL





QC: import/export processes

- When importing data into SQL from files or other databases, there is a significant risk of losing data through conversion errors
- It is a good practice to import a raw table with fields in a text format first, and create a SQL script to convert the data types in a new table to assist with reconciling to source data inputs



DATA SOCIETY © 2017

Note: Importing and then transforming data types allows simple reconciliation facilitates correction of conversion errors

INTRODUCTION TO SQL

Note: without a complete import there is no straight forward way to know if NULLs in a data set are true NULLs or conversion errors





Organization: commenting

- Syntax:

-- Selecting all contract numbers for comparison SELECT Contract NO FROM [tbl DS Sample]

query is run

SELECT Contract NO /* These contract numbers represent the population of non- performing loans */ FROM [tbl DS Sample]

INTRODUCTION TO SQL

 Code should always be commented to the extent that a non-SQL user can understand the analysis process, and for SQL users to repeat the process

- In-line comments – all text after "--" will not be read by SQL Server when a query is run

- Block comments - all text between "/*" and "*/" will not be read by SQL Server when a



Organization: formatting

Consistent code formatting makes it easier to read and edit. This saves hours of time spent trying to understand what is happening if the code is reviewed later. Here are a few suggestions for SQL formatting:

- Put keywords in UPPER CASE
 - SELECT, FROM
 - SUM, AVG
 - CASE, CAST, CONVERT
- Use tabs to align blocks of code
- Put commas at the beginning of a line

INTRODUCTION TO SQL

```
--Create Greater than 2000 flag
 SELECT A.Contract NO
  , CASE WHEN Origination Year > 2000 THEN `Y'
        ELSE 'N' END AS GT 2000 Flag
 INTO #TEST
 FROM [tbl DS Sample] A
 LEFT JOIN [tbl DS Sample2] B
   ON A.Contract NO = B. Contract NO
 WHERE A.Contract NO IS NOT NULL
--View records from #Test
  SELECT *
  FROM #TEST
```





Organization: naming conventions

- Be consistent when renaming fields within a table/project
 - If there are many analyses relying on a field in a table and the name changes, it could severely interrupt other processes before the proper changes are made across a team or organization
- Do not use spaces in names

around the name to be properly read in SQL

• Use consistent naming conventions for SQL tables, views, and other objects

- When spaces are used in a table or column names, it requires brackets ("[Field Name]")



Organization: saving SQL scripts

- Follow your office or team rules for creating and organizing scripts
- Each script should be modular and relate to one analysis or function
- Version control is also very important and adding version descriptions and dates to scripts can be very helpful when an analysis needs to be revisited
- For example a project named Surfside may have files saved like: - Surfside_Data_Prep_v01_20170225.sql = scripts related to preparing a table or data set for
 - analysis
 - Surfside Review v01_20170315.sql = scripts such as frequency distributions that describe the data
 - Surfside_Data_Analysis_v01_20170322.sql = scripts related to the project deliverable



Performance: argument order



INTRODUCTION TO SQL

- Try to use a leading character with LIKE (ex. LIKE 'm%' instead of LIKE '%m')
- Use LIKE instead of SUBSTRING with =

- Be careful with OR
- If multiple ANDs, put least likely condition first
- If equally likely, put least complex condition first



100

Performance: query structure

- Restrict result sets by using WHERE or only selecting the columns needed
- Use WHERE with HAVING when appropriate
- ORDER BY is inefficient; sort results in a separate step

INTRODUCTION TO SQL





Performance: space efficiency

- Backup your database, but don't store excessive copies of the backup
- Choose between temp tables and views
- Use appropriate field types (avoid using NVARCHAR)









INTRODUCTION TO SQL



Questions?



DATA SOCIETY®

Appendix: Bridging data between SQL and R



Bridging R and SQL

- R can be used to extract and write data to a SQL database
- central location



INTRODUCTION TO SQL

• This allows analysis results sets and interim data sets to easily accessed from a



ODBC Connection

- - Server.
- connection that R can use to query a database



SQL Con = odbcDriverConnect(paste('driver={SQL Server}) ,';server=W530THINKPAD\\SQLEXPRESS' ,';database=Data Society SQL CLass' < ,';trusted connection=true', sep = '')

INTRODUCTION TO SQL

 R can connect to SQL using an Open Database Connectivity (ODBC) connection - ODBC is a protocol that allows for a connection between data sources such as Microsoft SQL

• The RODBC package uses the odbcDriverConnect() function to create an ODBC



driver={SQL Server} indicates that we are connecting to a MS SQL server

The server you would like to connect to should be entered here after "server="

The database within the server that you would like to connect to should be entered here after "database="

true indicates that windows authentication will be used for server access permissions



Create a list of SQL tables

package, we can execute the SQL command SP tables which will store a dataframe containing all SQL tables and views in a database.

SQL Table List = sqlQuery(SQL Con ,'SP tables') 🗲

SQL Table List = SQL Table List[SQL Table List\$TABLE OWNER != 'sys' & SQL Table List\$TABLE OWNER != 'INFORMATION SCHEMA',]

INTRODUCTION TO SQL

• Using the ODBC connection and the sqlQuery() function from the RODBC

Execute and store query

SQL_Con is the ODBC connection previously defined SP_tables is a function in SQL server that returns all database tables and Views (system and otherwise)

Subset the data set

DATA SOCIETY © 2017

Remove System tables and Information schemas that would not typically be part of your analysis


Import SQL tables (individual)

package, we can execute a SELECT query which will store a dataframe containing all the results of that query.



INTRODUCTION TO SQL

• Using the ODBC connection and the sqlQuery() function from the RODBC

Execute and store SELECT queries

DATA SOCIETY © 2017



• Using a for loop the ODBC connection and the sqlQuery() function from the RODBC package, we can execute multiple SELECT queries to create dataframes for all (or a selected list) of tables and views in a database.

```
for(i in 1:nrow(SQL Table List)) {
schema = as.character(SQL Table_List$TABLE_OWNER[i])
tbl name = as.character(SQL Table List$TABLE NAME [i]
tbl import loop = sqlQuery(SQL Con
                           , paste('select * from
                                  , schema
                                  , tbl name
                           , sep = '')
assign(paste( schema,
             , tbl name , sep =
              ,tbl import loop
print(paste(schema, ' ', tbl name , sep = '')) +
```

Import SQL tables (bulk)

Import all tables



- get the schema name for each record to ensure the right object is being selected from the database
- get the table name for querying each table in the database for import
- code to import SQL table into R
- code to assign each table imported an appropriate variable name based on both the schema and the table name from the SQL database
- Print tables imported in the console to track progress while loop runs





Export data to SQL

• Using the ODBC connection and either the sqlSave() function or the sqlQuery() function from the RODBC package, we can save a data into a new SQL table or append the data to an existing data table from R.



INTRODUCTION TO SQL

Export dataframe to SQL table

- SQL connection created earlier in the code
- R dataframe being exported
- The name of the table being created and/or appended to
- excludes rownames from being exported as a column in the new SQL table
- TRUE if the table already exists then the records from this data frame will be appended to the existing table
- Define values
- SQL connection created earlier in the code
- Use INSERT statement with variable values to append data to a SQL table



110

Close ODBC connections

Close connections



INTRODUCTION TO SQL

• It is a best practice to close an ODBC connection (or connections) after you are done accessing the databases to help maintain the performance and availability of your SQL database. We use the functions odbcClose() and odbcCloseAll() from the RODBC package to close individual or all connections respectively.

> Close a specific ODBC connection Close all ODBC connection

DATA SOCIETY © 2017

111





INTRODUCTION TO SQL



Questions?

DATA SOCIETY © 2017

