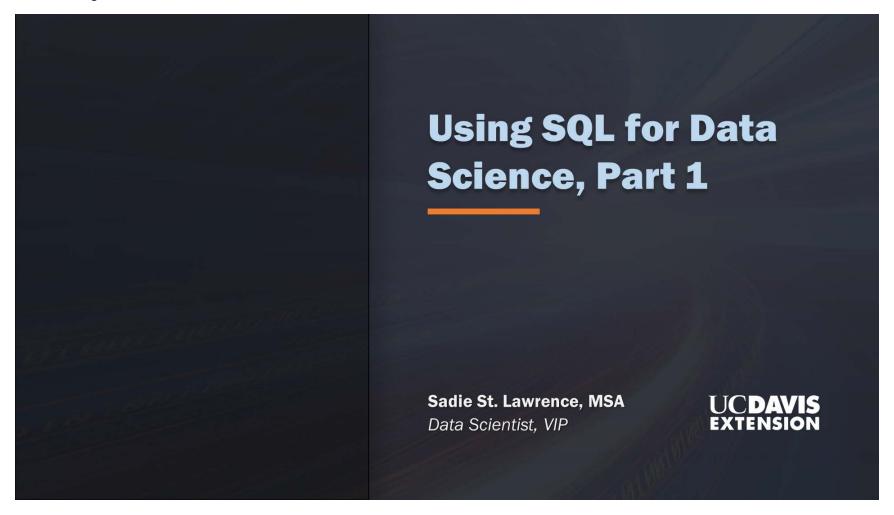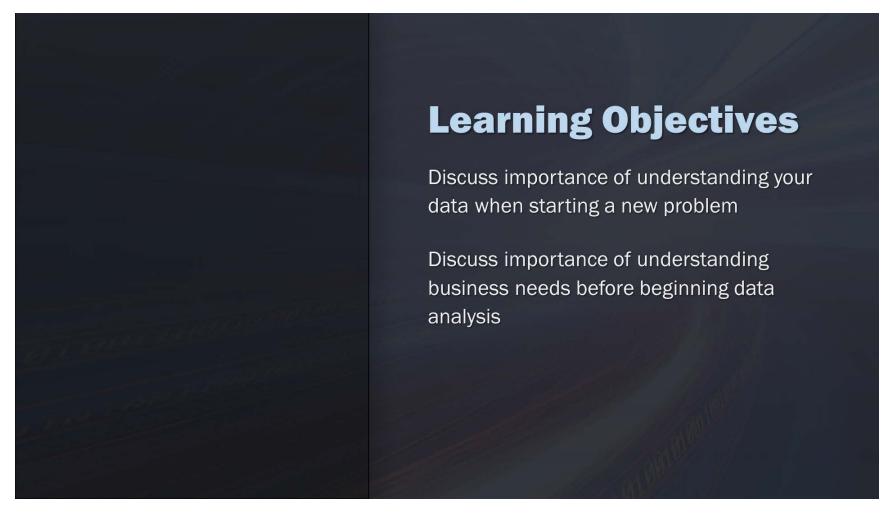Slide 1:  Using SQL for Data Science, Part 1

Slide 2: Learning Objectives
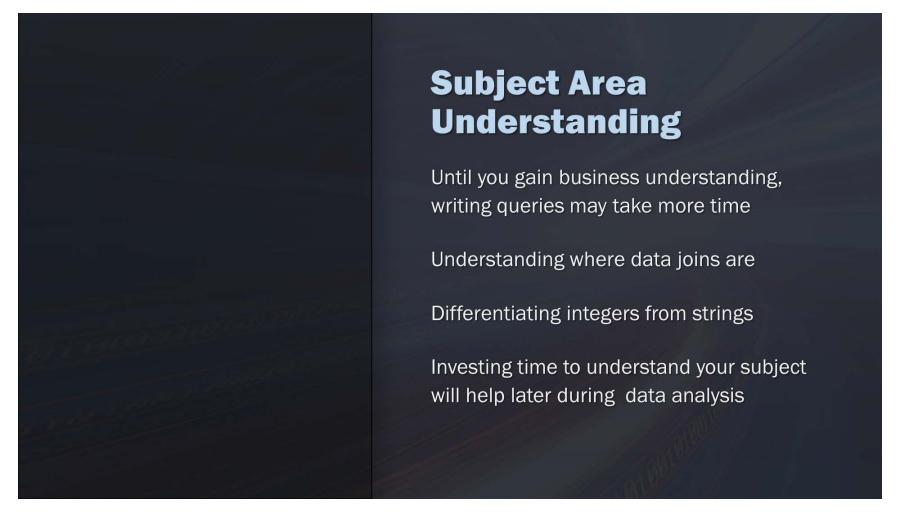
Slide 3:  Working Through a Problem from Beginning to End

Slide 4:  Data Understanding

Slide 5:  Subject Area Understanding



## Subject Area Understanding

Until you gain business understanding, writing queries may take more time

Understanding where data joins are

Differentiating integers from strings

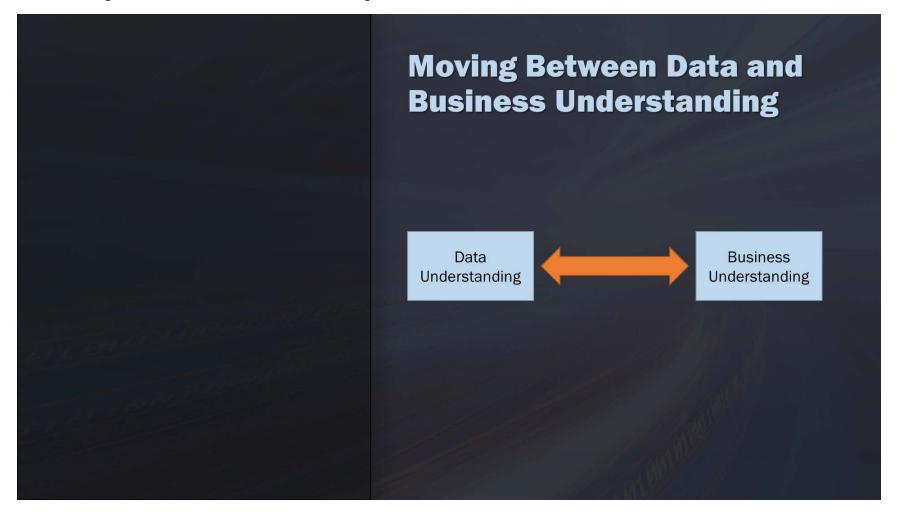Investing time to understand your subject will help later during  data analysis
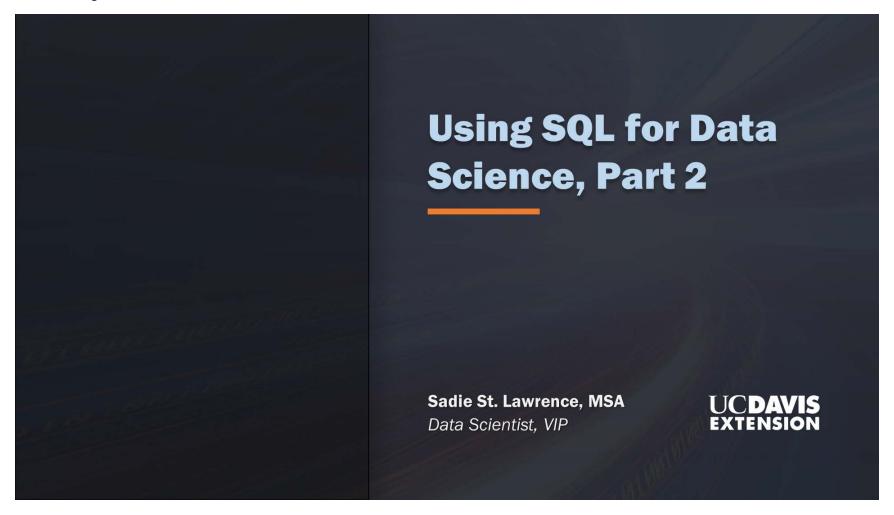
Slide 6:  Business Understanding

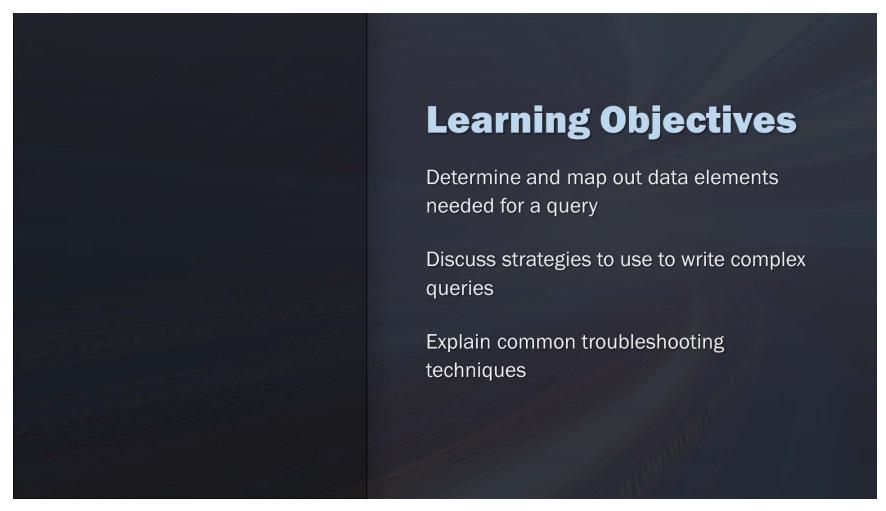Slide 7:  Beware of the Unspoken Need
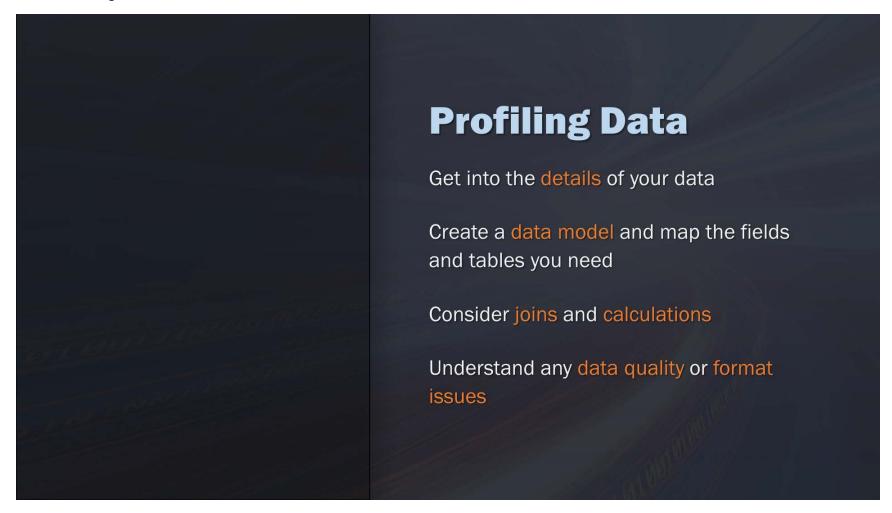
Slide 8:  Moving Between Data and Business Understanding

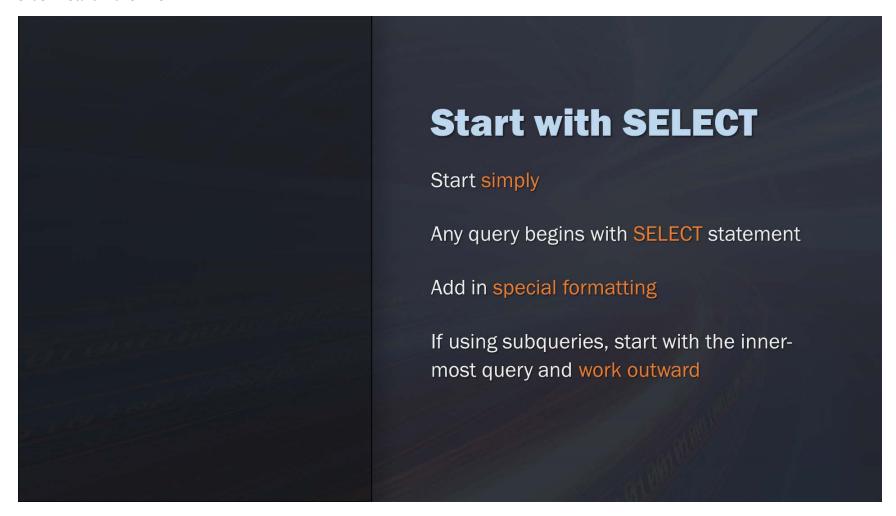Slide 1: Using SQL for Data Science, Part 2
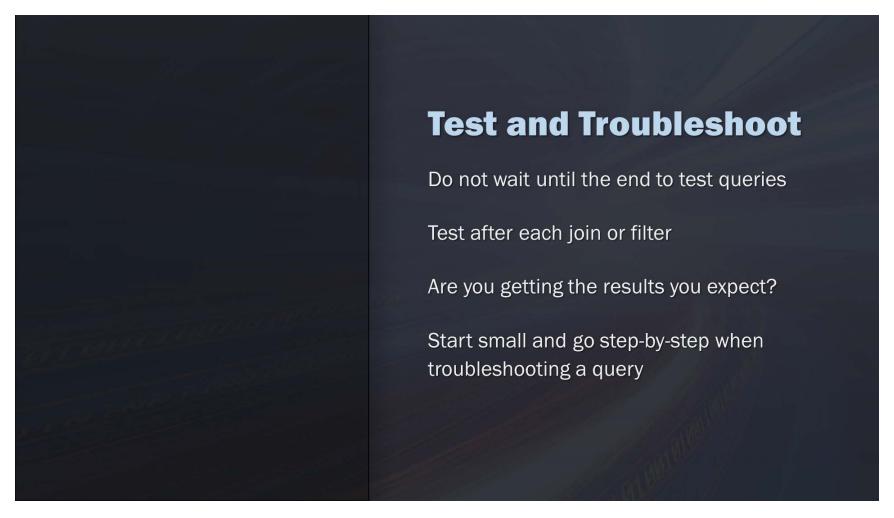
Slide 2:  Learning Objectives

Slide 3:  Profiling Data

Slide 4:  Start with SELECT

Slide 5:  Test and Troubleshoot

Slide 6:  Format and Comment

Slide 7: Review

Slide 1:  Data Governance and Profiling

Slide 2: Learning Objectives

Slide 3:  What is Data Profiling

Slide 4:  Object Data Profile

Slide 5: Column Data Profile

Slide 6:  Governance Best Practices



**Governance Best Practices**

Understand your read and write capabilities

Clean up your environments

Understand your promotion process

Slide 1:  Case Statements

Slide 2:  Learning Objectives

Slide 3:  What is a Case Statement

# What is a Case Statement

Mimics if-then-else statement found
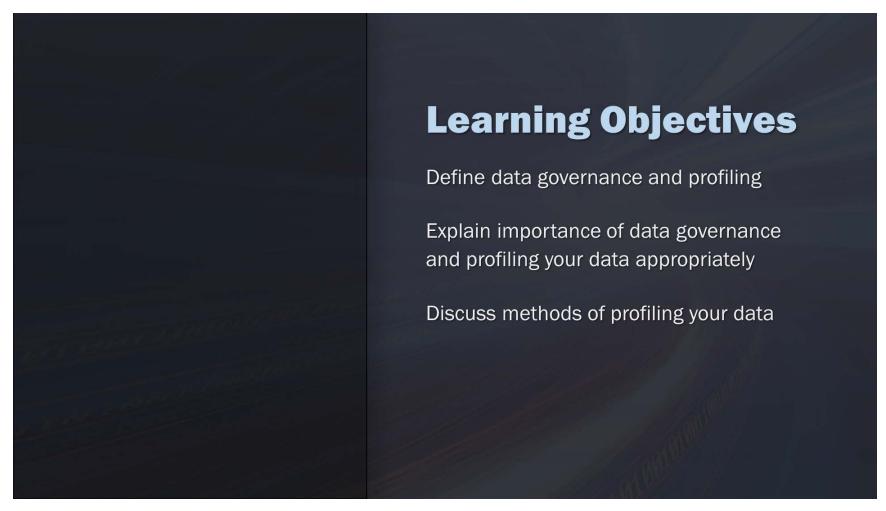in most programming languages

Can be used in SELECT, INSERT,
UPDATE, and DELETE statements

```
CASE
WHEN C1 THEN E1
WHEN C2 THEN E2
. . .
ELSE [result else]
END
```

```
CASE input_expression
        WHEN when_expression THEN result_expression [ ...n ]
        [ ELSE else_result_expression ]
END
```

Slide 4: Simple Case Statement

Slide 5: Search Case Statement

Slide 6:  Search Case Statement

Slide 1: Working with Date and Time Strings

Slide 2: Learning Objectives



# Learning Objectives

Describe the complexities of adjusting date and time strings

Discuss the different formats in which dates and times are presented

List and describe the 5 different functions in SQL that can be used to manipulate date and time strings

Slide 3: Working with Date Variables

# Working with Date Variables

"As long as your data contains only the date portion, your queries will work as expected. However, if a time portion is involved, it gets more complicated."

"The most difficult part when working with dates is to be sure that the format of the date you are trying to insert, matches the format of the date column in the database."

-W3 Schools

Dates are stored as datetypes

Each DBMS uses it's own variety of datatypes

```
Wednesday, September 17th, 2008
9/17/2008 5:14:56 P.M. EST
9/17/2008 19:14:56 GMT
2612008 (Julian format)
```

Slide 4:  Date Formats



Date Formats

DATE
Format YYYY-MM-DD

DATETIME
Format: YYYY-MM-DD HH:MI:SS

TIMESTAMP
Format: YYYY-MM-DD HH:MI:SS

If you query a DATETIME with:

```
WHERE PurchaseDate='2016-12-12'
```

You will get no results

Slide 5:  SQLite Date Time Functions

Slide 6:  Timestrings

Slide 7:  Modifiers



Modifiers

| | |
|---|---|
| NNN days | start of year |
| NNN hours | start of day |
| NNN minutes | weekday N |
| NNN.NNNN seconds | unixepoch |
| NNN months | localtime |
| NNN years | utc |
| start of month | |

Slide 1: Views

Slide 2:  Learning Objectives

Slide 3:  Overview of Views

# Overview of Views

A stored query

Can add or remove columns without changing schema

Use it to encapsulate queries

The view will be removed after database connection has ended

```
CREATE [TEMP] VIEW [IF NOT EXISTS]
view_name(column-name-list)
AS
select-statement;
```

Slide 4:  Creating a View

Slide 5: Creating a View

Slide 6:  Why Use Views



# Why Use Views

Get a count of how many territories each employee has

| | count(territorydescription) | Lastname | Firstname |
|---|---|---|---|
| 1 | 7 | Buchanan | Steven |
| 2 | 4 | Callahan | Laura |
| 3 | 2 | Davolio | Nancy |
| 4 | 7 | Dodsworth | Anne |
| 5 | 7 | Fuller | Andrew |
| 6 | 10 | King | Robert |
| 7 | 4 | Leverling | Janet |
| 8 | 3 | Peacock | Margaret |
| 9 | 5 | Suyama | Michael |

```
SELECT count(territorydescription)
,Lastname
,Firstname
FROM my_view
GROUP BY Lastname, Firstname;
```

Slide 1: Date and Time String Examples

Slide 2:  Learning Objectives



# Learning Objectives

Use the STRFTIME function

Compute current date and use it to compare to a recorded date in your data

Use the NOW function

Combine several date and time functions together to manipulate data
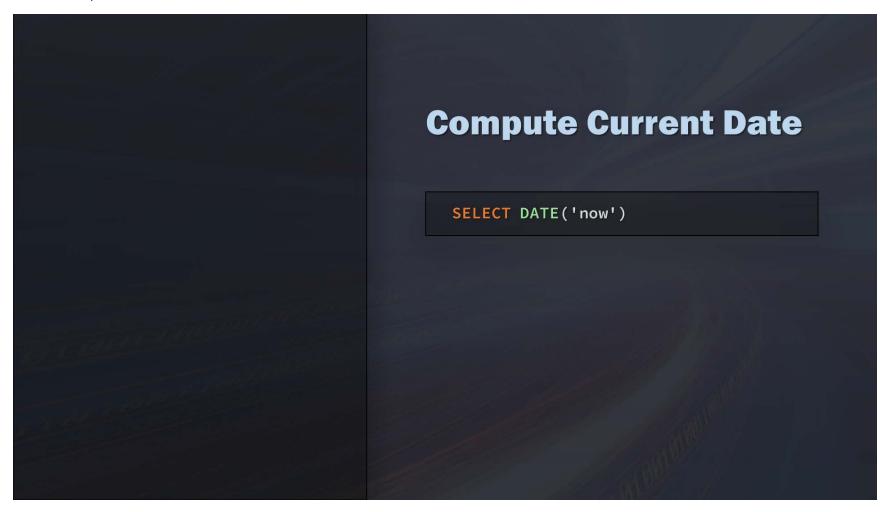
Slide 3:  Example

Slide 4:  Compute Current Date

Slide 5:  Compute Year, Month, and Date for the Current Date



**Compute Year, Month and Day for the Current Date**

```
SELECT STRFTIME('%Y %m %d','now')
```

Slide 6:  Compute the Hour, Minute and Second and Milliseconds from Current DATETIME

Slide 7:  Compute Age Using Birthdate