



DATABASES AND SQL WITH PANDAS

A "SUPER PYTHON" TALK BY NICHOLAS A. DEL GROSSO

TODAY'S GOALS

1. Understand the Benefits of Tables as Data Structures
 - ☐ 5-minute Geek Out Session
 - ☐ 5-minutes of Pulling Myself Together
2. Review Pandas as a Tool for working with Tabular data.
 - ☐ 20-minute guided, hands-on exercise
3. Become Familiar with the Concept of a "Relational" Database.
 - ☐ Demonstration of Normalization
4. Understand where Database software fits into all this.
 - ☐ 10-Minute Live Coding Demo: SQLite3
5. Learn some SQL and how to use it with Pandas.
 - ☐ 10-minute Live Coding Demo
 - ☐ 30-minute Exercises

TABLE: EXPERIMENT

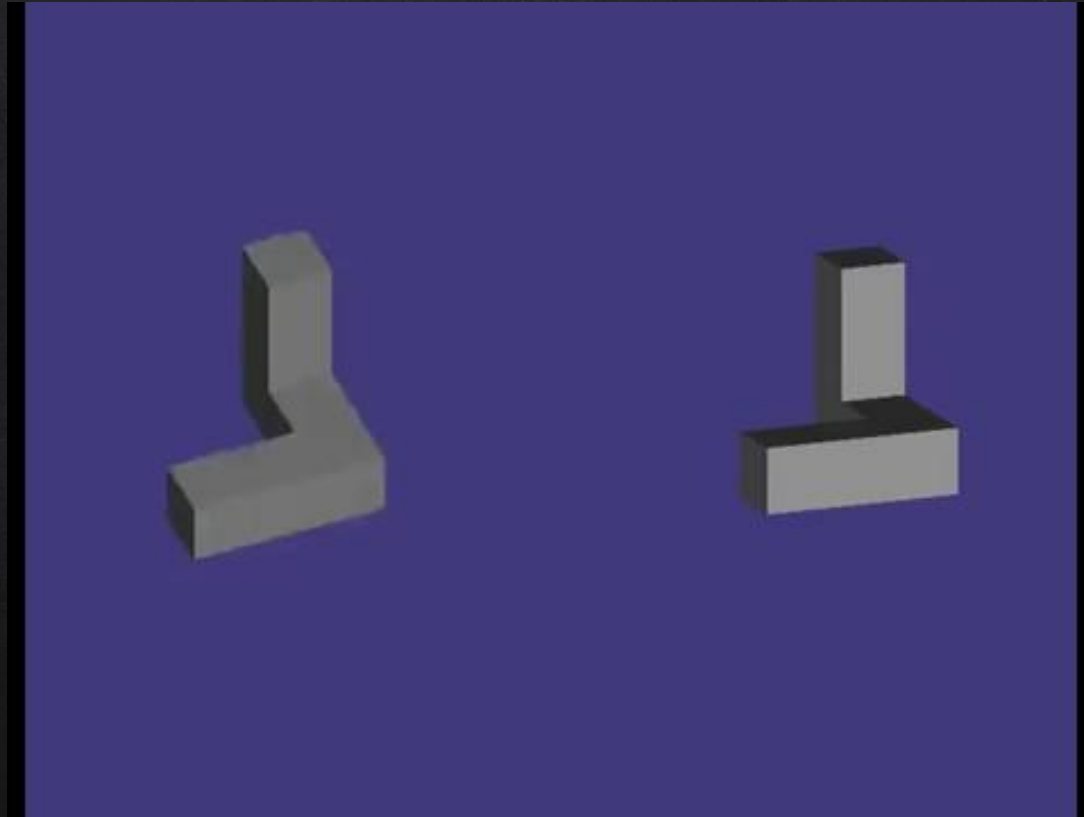
Subject	Gender	Session	Date	Trial	Condition	Drink	Reaction Time
Nick	M	1	5/5/16	1	Exp	Beer	100.2
Nick	M	1	5/5/16	2	Ctrl	Water	83.4
Nick	M	1	5/5/16	3	Ctrl	Water	95.2
Nick	M	1	5/5/16	4	Exp	Beer	78.0
Nick	M	1	5/5/16	5	Exp	Water	104.3
Nick	M	1	5/5/16	6	Ctrl	Water	87.9

ALGORITHMIC CONCEPTS: QUERYING A TABLE

- ❑ How many total Trials were there, across the entire Experiment?
- ❑ How many Subjects are in the Experiment?
- ❑ What was the mean Reaction Time?
 - ...for each Condition?
 - ...for each Subject?
 -for each Subject, for each Condition?

Subject	Gender	Session	Date	Trial	Condition	Stimulus	Reaction Time
Nick	M	1	5/5/16	1	Exp	Beer	100.2
Nick	M	1	5/5/16	2	Ctrl	Water	83.4

PANDAS REVIEW: LOADING AND QUERYING TABULAR DATA



PROBLEMS WITH THE TABLE: REDUNDANCY AND LARGE FILE SIZE

Solution: “Normalization” and the “Relational Database”

EXPERIMENT

Subject	Gender	Session	Date	Trial	Condition	Stimulus	Reaction Time
Nick	M	1	5/5/16	1	Exp	Beer	100.2
Nick	M	1	5/5/16	2	Ctrl	Water	83.4
Nick	M	1	5/5/16	3	Ctrl	Water	95.2
Nick	M	1	5/6/16	4	Exp	Bear	78.0
Nick	M	1	5/5/16	5	Exp	Water	104.3
Nick	M	1	5/5/16	6	Ctrl		87.9

SUBJECTS

Name	Gender
Nick	M
Anna	F

CONDITIONS

Condition	Drink
Exp	Beer
Ctrl	Water

SESSIONS

Session	Date
1	5/5/16
1	5/5/16

TRIALS

Trial	Reaction Time
1	100.2
2	83.4
3	95.2
4	78.0
5	104.3
6	87.9

SUBJECTS

ID	Name	Gender
1	Nick	M
2	Anna	F

CONDITIONS

ID	Condition	Drink
1	Exp	Beer
2	Ctrl	Water

SESSIONS

ID	Subject	Session	Date
1	1	1	5/5/16
2	2	1	5/5/16

TRIALS

ID	Session	Condition	Trial	Reaction Time
1	1	1	1	100.2
2	1	2	2	83.4
3	1	2	3	95.2
4	1	1	4	78.0
5	1	1	5	104.3
6	1	2	6	87.9

QUERYING A RELATIONAL DATABASE: MAKE A SINGLE TABLE FIRST THROUGH THE "JOIN"

- "WHICH PARTICIPANT WAS IN SESSION 1?"
 - JOIN SESSIONS AND SUBJECTS,
 - TAKE ONLY NAME AND SESSION

Session	Name
1	Nick
2	Anna

- "WHICH STIMULI WERE USED IN EACH SESSION?"
 - JOIN CONDITIONS AND SESSIONS AND TRIALS,
 - TAKE ONLY SESSION AND STIMULUS

Session	Drink
1	Beer
1	Water
2	Beer

QUICK DEMO: JOIN IN PANDAS

```
pd.merge(conditions, trials, left_on = 'ID', right_on =  
'Condition')
```

```
trials.join(conditions, on='Condition')
```

REVIEW SUMMARY: NORMALIZATION

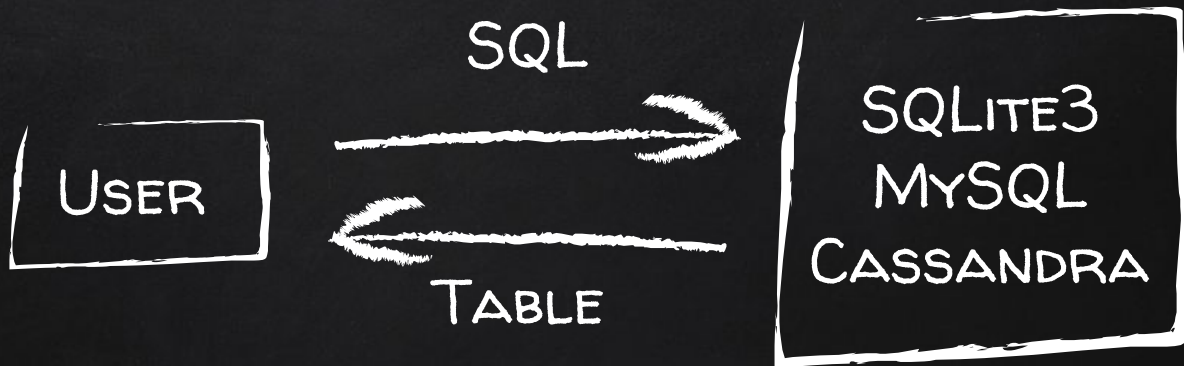
- Normalization for Storage: 1 Table → Multiple Tables
 - Reduces File Size
 - Reduces Redundancy
 - Reduces Errors
 - Increases Read / Write Speeds
 - Can Form a More “Natural” Organization Schema
- Querying through Joins: Multiple Tables → 1 Table

SQL



Apache
CASSANDRA™

SQL IS THE LANGUAGE USED TO TALK TO RELATIONAL
DATABASE SOFTWARE



SQL: STRUCTURED QUERY LANGUAGE

Create Tables

- CREATE TABLE
Define a new table
- DROP TABLE
Remove a table

Query and Update Columns

- SELECT
Retrieve columns from a table or view
- INSERT INTO
Create new rows in a table
- COMMIT
Save changes to the database
- UPDATE
Update rows of a table

Filter, Group, and Process the Rows

- WHERE
To retrieve specific information from a table excluding other irrelevant data
- DISTINCT
To return only distinct (different) values in a column
- BETWEEN..AND..
To select values within a range

SQL QUERYING EXAMPLES

```
SELECT Name, Gender FROM Subjects;
```

Name	Gender
Nick	M
Anna	F

```
SELECT * FROM Conditions;
```

ID	Condition	Drink
1	Exp	Beer
2	Ctrl	Water

```
SELECT ReactionTime FROM Trials LIMIT 2;
```

Reaction Time
100.2
83.4

SQL QUERYING EXAMPLES:

FILTERING, AGGREGATING, AND JOINING

```
SELECT Name FROM Subjects  
WHERE Gender = "M";
```

Name
Nick

```
SELECT Condition, avg(ReactionTime) FROM Trials  
GROUP BY Condition;
```

Condition	avg(ReactionTime)
1	78.486
2	92.112

```
SELECT Drink, avg(ReactionTime) FROM Trials  
JOIN Trials ON Conditions.ID = Trials.Condition  
GROUP BY Condition;
```

Drink	avg(RT)
Beer	78.486
Water	92.112

A COUPLE SQLITE-SPECIFIC COMMANDS

Get all Table Names:

```
SELECT name FROM SQLITE_MASTER;
```

Get Column Names from a Table:

```
PRAGMA table_info(<TableName>);
```

DEMO: SQL QUERYING IN PYTHON

Pandas uses SQLAlchemy to connect to Databases like SQLite3

```
from sqlalchemy import create_engine
engine = create_engine('sqlite:///my_folder/my_data.db')
conn = engine.connect()
```

```
import pandas as pd
query = "SELECT Angle, Correct, Matching FROM Trials;"
df = pd.read_sql(query, conn)
```

```
import seaborn as sns
sns.factorplot(x='Angle', y='Correct', hue='Matching',
data=df)
```


30 MINUTES HANDS-ON: (THANK YOU FOR YOUR ATTENTION!)

- ☐ How Many Subjects are in this Study?
- ☐ Did Subjects Take Longer to Respond to Matching or Nonmatching Stimuli?
- ☐ Was there a ReactionTime – Stimulus Rotation Relationship?
- ☐ What was the mean Subject Age?
- ☐ Was there an Effect of Subject Sex on:
 - Reaction Time
 - Accuracy (Correctness)