

DATABASES

Sri Kanajan

DATABASES

LEARNING OBJECTIVES

- Understanding of the uses and differences of databases at a high level
- Accessing databases from Pandas
- Understanding the basics of SQL

OPENING

DATABASES

DATABASES

- Today's lesson will be on databases and the SQL query language.
- Databases are the standard solution for data storage. They're far more robust than text and CSV files.
- They come in many flavors, but we'll explore the most common: *relational databases*.



DATABASES

- Relational databases also come in different varieties, but almost all use SQL as a basis for querying (i.e. retrieving) data.
- Most analyses typically involve pulling data from a database.

INTRODUCTION

DATABASES

DATABASES

- Databases are computer systems that manage the storage and querying of datasets.
- They provide a way to organize the data on disk (i.e. hard drive) and efficient methods to retrieve information. Databases allow a user to create rules that ensure proper data management and verification.
- Typically, retrieval is performed using a query language, a mini programming language with a few basic operators for data transformation.
- The most common query language is SQL (Structured Query Language).

DATABASES

- A *relational database* is based on links between data entities or concepts.
- Typically, a relational databases is organized into *tables*.
- Each table should correspond to one entity or concept. Each table is similar to a single CSV file or Pandas dataframe.
- For example, consider an application like Twitter. Our two main entities are Users and Tweets. For each of these, we would have a separate table.

DATABASES

- A table is made up of rows and columns, similar to a Pandas dataframe or Excel spreadsheet.
- Each table has a specific *schema*, a set of rules for what goes in each table. These specify which columns are contained in the table and what *type* of data is in each column (e.g. text, integers, decimals, etc).

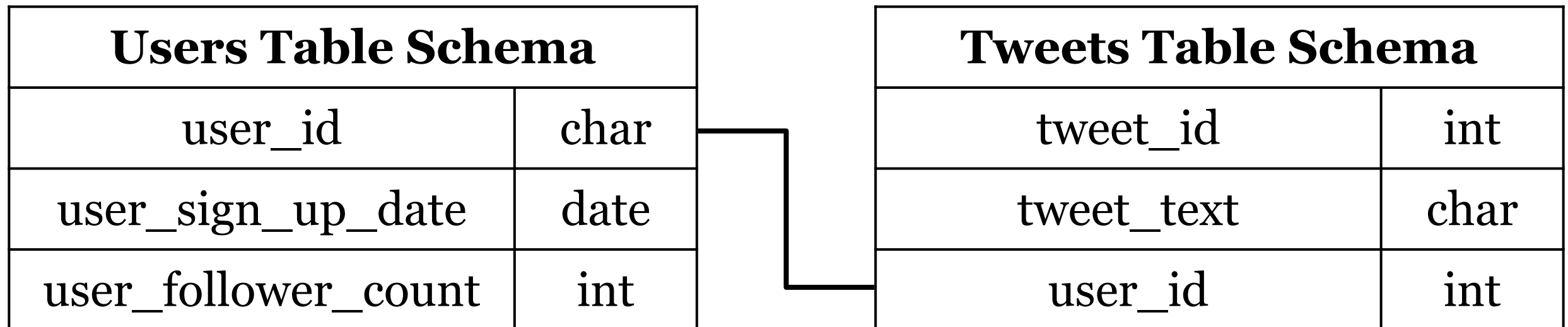
Users Table Schema	
user_id	char
user_sign_up_date	date
user_follower_count	int

DATABASES

- Each table typically has a *primary* key column. This column has a unique value per row and serves as the identifier for the row.
- A table can have many *foreign keys* as well. A *foreign key* is a column that contains values to link the table to the other tables.
- These keys that link the table together define the relational database.

DATABASES

- For example, the tweets table may have as columns:
 - tweet_id - the primary key tweet identifier
 - tweet_text
 - user_id - a foreign key to the users table



DATABASES

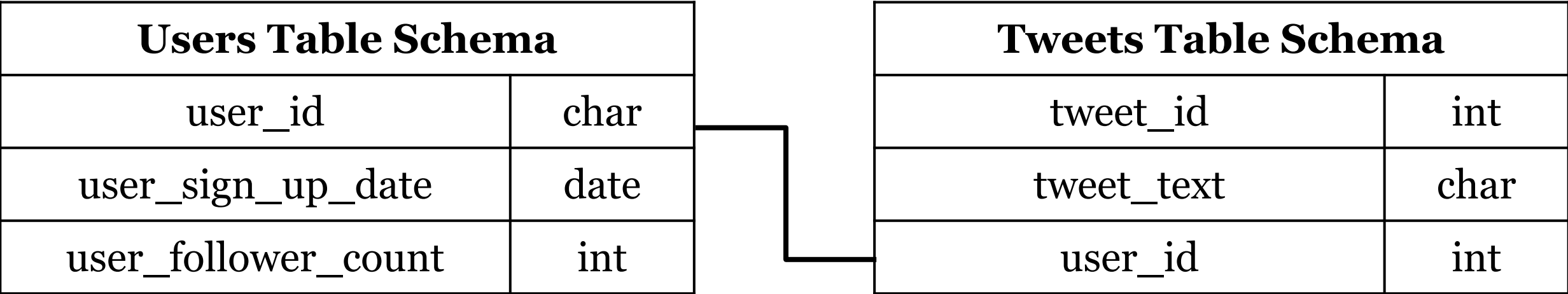
- MySQL and Postgres are popular variants of relational databases and are widely used. Both are open-source and available for free.
- Alternatively, many companies use proprietary software such as Oracle or Microsoft SQL databases.
- While these databases offer many of the same features and use the same SQL language, the latter two offer some maintenance features and support that large companies find useful.

NORMALIZED VS DENORMALIZED DATA

- Once we start organizing our data into tables, we start to separate it into *normalized* and *denormalized* setups.
- *Normalized* structures have a single table per entity and use many foreign keys or link tables to connect the entities.
- *Denormalized* structures have fewer tables that combine different entities.

NORMALIZED VS DENORMALIZED DATA

- With our Twitter example, a *normalized* structure would place users and tweets in different tables.



NORMALIZED VS DENORMALIZED DATA

▸ A *denormalized* structure would put them both in one table.

Twitter Table Schema	
tweet_id	int
tweet_text	char
user_id	int
user_sign_up_date	date
user_follower_count	int

NORMALIZED VS DENORMALIZED DATA

Denormalized structures:

- Duplicates a lot of information
- Makes data easy to access since it's all in one table

Normalized structures:

- Save storage space by separating information
- Requires joining of table to access information about two different entities, a slow operation

DEMO

ACCESSING DATABASES FROM PANDAS

ACCESSING DATABASES FROM PANDAS

- While databases provide many analytical capabilities, often it's useful to pull the data back into Python for more flexible programming.
- Large, fixed operations would be more efficient in a database, but Pandas allows for interactive processing.
- This would run very efficiently in a database vs connecting to Python.

ACCESSING DATABASES FROM PANDAS

- However, if we want to investigate the login or sales data further and ask more interactive questions, then using Python would come in very handy.

```
import pandas as pd  
from pandas.io import sql
```

- Pandas can be used to connect to most relational databases.

ACCESSING DATABASES FROM PANDAS

- In this demonstration, we will create and connect to a SQLite database. SQLite creates portable relational databases saved in a single file.
- These databases are stored in a very efficient manner and allow fast querying, making them ideal for small databases or databases that need to be moved across machines.
- Additionally, SQLite databases can be created with the setup of MySQL or Postgres databases.

WRITING DATA INTO A DATABASE

- Data is moved to the database with the `to_sql` command, similar to the `to_csv` command.
- `to_sql` takes several arguments.
 - `name` - the table name to create
 - `con` - a connection to a database
 - `index` - whether to input the index column
 - `schema` - if we want to write a custom schema for the new table
 - `if_exists` - what to do if the table already exists. We can overwrite it, add to it, or fail

ACTIVITY: KNOWLEDGE CHECK

ANSWER THE FOLLOWING QUESTIONS



EXERCISE

1. Load the Rossmann Store metadata in `rossmann-stores.csv` and create a table in the database with it.

DELIVERABLE

Created table for store metadata

DEMO

**SQL SYNTAX: SELECT,
WHERE, GROUP BY,
JOIN**

SQL OPERATORS: SELECT

- Every query should start with SELECT. SELECT is followed by the names of the columns in the output.
- SELECT is always paired with FROM, which identifies the table to retrieve data from.

```
SELECT <columns>  
FROM <table>
```

- SELECT * denotes returning *all* of the columns.

SQL OPERATORS: SELECT

▸ Rossmann Stores example:

```
SELECT Store, Sales  
FROM rossmann_sales;
```

SQL OPERATORS: WHERE

- WHERE is used to filter a table using a specific criteria. The WHERE clause follows the FROM clause.

```
SELECT <columns>  
FROM <table>  
WHERE <condition>
```

- The condition is some filter applied to the rows, where rows that match the condition will be output.

SQL OPERATORS: WHERE

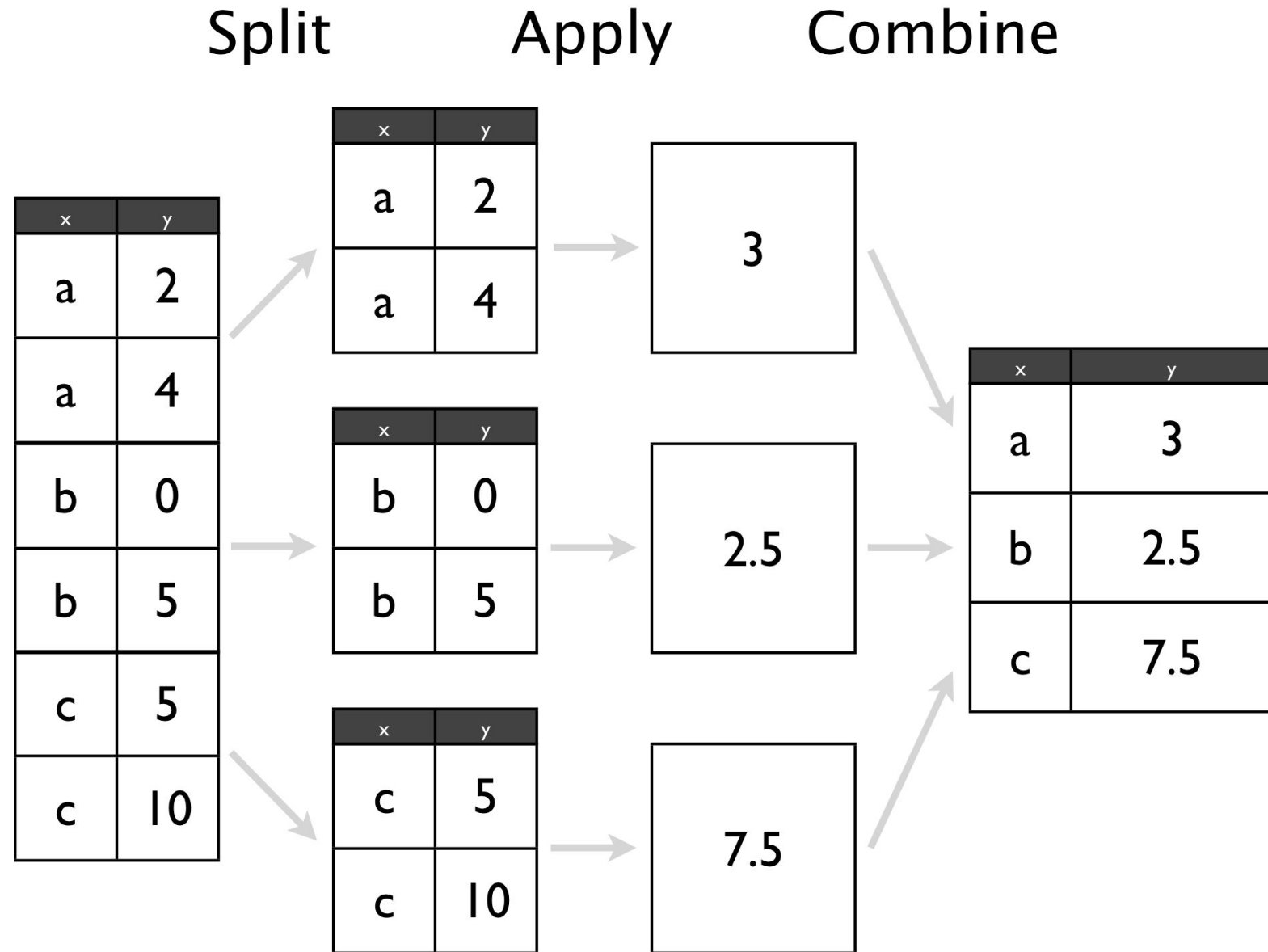
▸ Rossmann Stores example:

```
SELECT Store, Sales  
FROM rossmann_sales  
WHERE Store = 1;
```

```
SELECT Store, Sales  
FROM rossmann_sales  
WHERE Store = 1 and Open = 1;
```

SQL OPERATORS: GROUP BY

- ▶ GROUP BY allows us to aggregate over any field in the table by applying the concept of Split Apply Combine.
- ▶ We identify some key with which we want to segment the rows. Then, we roll up or compute some statistics over all of the rows that match that key.



SQL OPERATORS: GROUP BY

- GROUP BY *must* be paired with an aggregate function, the statistic we want to compute in the rows, in the SELECT statement.
- COUNT(*) denotes counting up all of the rows. Other aggregate functions commonly available are AVG (average), MAX, MIN, and SUM.
- If we want to aggregate over the entire table, without results specific to any key, we can use an aggregate function in the SELECT clause and ignore the GROUP BY clause.

SQL OPERATORS: GROUP BY

▸ Rossmann Stores example:

```
SELECT Store, SUM(Sales), AVG(Customers)
FROM rossmann_sales
WHERE Open = 1
GROUP BY Store;
```

SQL OPERATORS: ORDER BY

- ORDER BY is used to sort the results of a query.

```
SELECT <columns>  
FROM <table>  
WHERE <condition>  
ORDER BY <columns>
```

- You can order by multiple columns in ascending (ASC) or descending (DESC) order.

SQL OPERATORS: ORDER BY

- Rossmann Stores example:

```
SELECT Store, SUM(Sales) as total_sales, AVG(Customers)
FROM rossmann_sales
GROUP BY Store
WHERE Open = 1;
ORDER BY total_sales desc;
```

- COUNT(*) AS cnt renames the COUNT(*) value to cnt so we can refer to it later in the ORDER BY clause.

SQL OPERATORS: JOIN

- JOIN allows us to access data across many tables. We specify how a row in one table links to another.

```
SELECT a.Store, a.Sales, s.CompetitionDistance  
FROM rossmann_sales a  
JOIN rossmann_stores s  
ON a.Store = s.Store
```

- Here, ON denotes an *inner* join.

SQL OPERATORS: JOIN

- By default, most joins are an *Inner Join*, which means only when there is a match in both tables does a row appear in the results.
- If we want to keep the rows of one table *even if there is no matching counterpart*, we can perform an *Outer Join*.
- Outer joins can be LEFT, RIGHT, or FULL, meaning keep all of the left rows, all the right rows, or all the rows, respectively.

SQL OPERATORS: WINDOW FUNCTIONS

- A way to do a group by within the select clause itself and still maintain the original primary key
- `SELECT depname, empno, salary, avg(salary) OVER (PARTITION BY depname) FROM empsalary;`

depname	empno	salary	avg
develop	11	5200	5020.0000000000000000
develop	7	4200	5020.0000000000000000
develop	9	4500	5020.0000000000000000
develop	8	6000	5020.0000000000000000
develop	10	5200	5020.0000000000000000
personnel	5	3500	3700.0000000000000000
personnel	2	3900	3700.0000000000000000
sales	3	4800	4866.666666666666667
sales	1	5000	4866.666666666666667
sales	4	4800	4866.666666666666667

(10 rows)

SQL OPERATORS: WINDOW FUNCTIONS

```
SELECT depname, empno, salary, enroll_date  
FROM  
  (SELECT depname, empno, salary, enroll_date,  
    rank() OVER (PARTITION BY depname ORDER BY salary DESC, empno) AS pos  
  FROM empsalary  
  ) AS ss  
WHERE pos < 3;
```

- What does this query do?

ADVANCED SQL

- Other functions such as LAG, LEAD help in time series analysis.
- Query planning
- Query performance analysis
- Indexing
- UDF - user defined queries

ACTIVITY: KNOWLEDGE CHECK

ANSWER THE FOLLOWING QUESTIONS



EXERCISE

1. Write a query for the Rossmann Sales data that returns Store, Date, and Customers.

DELIVERABLE

The requested query

ACTIVITY: KNOWLEDGE CHECK

ANSWER THE FOLLOWING QUESTIONS



EXERCISE

1. Write a query for the Rossmann Sales data that returns Store, Date, and Customers for stores that were open and running a promotion.

DELIVERABLE

The requested query

ACTIVITY: KNOWLEDGE CHECK

ANSWER THE FOLLOWING QUESTIONS



EXERCISE

1. Write a query that returns the total sales on the promotion and non-promotion days.

DELIVERABLE

The requested query

INDEPENDENT PRACTICE

PANDAS AND SQL

ACTIVITY: PANDAS AND SQL



EXERCISE

DIRECTIONS (40 minutes)

1. Load the Walmart sales and store features data.
2. Create a table for each of those datasets.
3. Select the store, date and fuel price on days it was over 90 degrees.
4. Select the store, date and weekly sales and temperature.
5. What were average sales on holiday vs. non-holiday sales?
6. What were average sales on holiday vs. non-holiday sales when the temperature was below 32 degrees?

DELIVERABLE

Answers to the above questions

CONCLUSION

TOPIC REVIEW

CONCLUSION

- While this was a brief introduction, databases are often at the core of any data analysis. Most analysis starts with retrieving data from a database.
- SQL is a key language that any data scientist should understand.
 - SELECT: Used in every query to define the resulting columns
 - WHERE: Filters rows based on a given condition
 - GROUP BY: Groups rows for aggregation
 - JOIN: Combines two tables based upon a given condition

CONCLUSION

- Pandas can be used to access data from databases as well. The result of the queries will end up in a Pandas dataframe.
- There is much more to learn about query optimization if one dives further!

LESSON

Q & A

LESSON

EXIT TICKET

DON'T FORGET TO FILL OUT YOUR EXIT TICKET