# Databases, Scrapping, and APIs

*Ivan Corneillet*

*Data Scientist*

# Learning Objectives

After this lesson, you should be able to:

- ‣ Understand the uses and differences of the major types of databases, including RDBMS databases (i.e., SQL) and NoSQL databases

- ‣ Access databases from *pandas*

- ‣ Describe how web scraping works, conceptually and explain practically how to web scraping works using Python

- ‣ Define how to approach scraping project data

- ‣ Describe APIs and how to make calls and consume API data
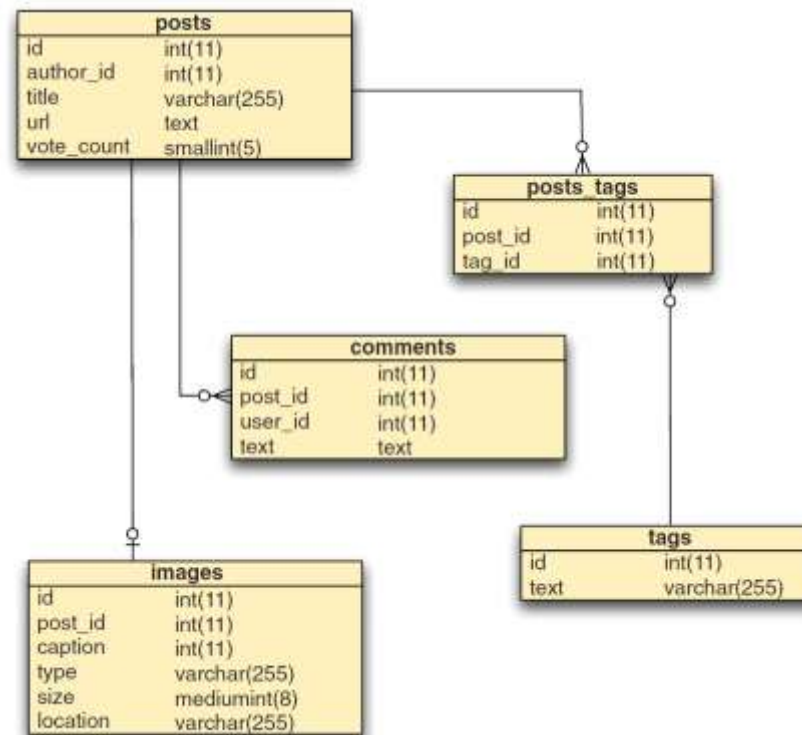
# Here's what's happening today:

- RDBMS/SQL Databases

  - The `SELECT` statement

- NoSQL (Not-Only SQL) Databases

  - NoSQL Classification

  - CAP Theorem

  - ACID vs. BASE

- Popular NoSQL Databases

- CRUD and REST

- Map Reduce

- Scrapping

- Application Programming Interfaces (APIs)

# Relational Database Management Systems (RDBMS)
# and
# Structured Query Language (SQL)

# A relational database links data entities and concepts. E.g., a relational data model for entries on a social news site



Source: MongoDB in Action

# Relational databases are organized into *tables.* Each table has a specific *schema,* a set of rules for what goes in each table

- Each table corresponding to one entity or concept

- A table is made up of rows and columns, similar to a *pandas* dataframe

- Schemas specify which columns are contained in the table and what type of data is in each column (e.g., text, integer, or date). This means you can't add text data to an integer column

- For this reason and many others, databases allow for stronger consistency of the data and are often a better solution for data storage

# CRUD and SQL

| CRUD(*)<br>(*) the four basic functions of persistent storage | SQL Statement | |
|:---:|:---:|:---:|
| Create | CREATE | Create or add new entries |
| Read | **SELECT** | Read, retrieve, search, or view existing entries |
| Update | UPDATE | Update or edit existing entries |
| Delete | DELETE | Delete/deactivate/remove existing entries |

# SQLite's **SELECT** Statement

(https://www.sqlite.org/lang_select.html)

# NoSQL (Not-Only SQL) Databases

# NoSQL databases fall into four primary classifications

‣ **Key-value stores** – use a simple data model that pairs a unique key and its associated value in storing data elements.  Common uses include storing clickstream data and application logs

‣ **(Wide)-column stores** (a.k.a., table-style databases) – store data across tables that can have very large numbers of columns.  Common uses include Internet search and other large-scale Web applications

‣ **Document databases** – store data elements in document-like structures that encode information in formats such as JSON.  Common uses include content management and monitoring Web and mobile applications

‣ **Graph databases** – emphasize connections between data elements, storing related "nodes" in graphs to accelerate querying.  common uses include recommendation engines and geospatial applications

# The CAP theorem states that in the presence of a network partition (P), one has to choose between consistency (C) and availability (A)



**Total Redundancy**
Systems stays up even when nodes fail
(i.e., the system gives a reasonable response in a reasonable time)

Availability

A

CA

AP

Pick Two!

C

CP

P

Consistency

Partition Tolerance

**ACID Transactions**
Everyone always have the same view of the data (regardless of any pending updates or deletes)

**Infinite Scaleout**
System stays up even when network between nodes fail (i.e., only a total network failure can cause the system to respond incorrectly)

# ACID vs. BASE: The pH of Database Transaction Processing

## ACID

‣ **Atomicity** – all operations are performed or none of them are. If one part of the transaction fails, then all fail

‣ **Consistency** – a transaction must meet all rules defined by the system at all times; there are never any half-completed transactions

‣ **Isolation** – transactions are independent from each other

‣ **Durability** – once complete, a transaction cannot be undone

## BASE

‣ **Basically Available** – the system will give and accept queries and give responses even in regards to node failures

‣ **Soft State** – the data is in a constant state of flux and might be stale

‣ **Eventual Consistency** – the data will eventually be consistent through all nodes and in all databases, but not every transaction at every moment

# AC/AP/CP vs. ACID/BASE

### AC

- Small datasets can be both consistent and available but a non-option in distributed systems (networks aren't completely reliable so you must tolerate partitions)

### AP
### BASE w/ eventual consistency

- System returns the most recent version of the data it has (which could be stale). The system will also accept writes that can be processed later when the partition is resolved
- Choose AP (over C) when you are flexible on when the data in the system synchronizes

### CP
### ACID w/ eventual availability

- System waits for a response from the partitioned node which could result in a timeout error
- Choose CP (over A) when you require atomic reads and writes

# Popular NoSQL databases<sup>(*)</sup>

| | | |
|---|---|---|
| Accumulo | Wide-column store | CP/ACID |
| Cassandra | Wide-column store | AP/BASE |
| CouchDB | Document database | AP/BASE |
| DynamoDB | Key-value store | AP/BASE |
| HBase | Wide-column store | CP/ACID |
| MongoDB | Document database | CP/ACID |
| Neo4j | Graph database | CP/ACID |
| Redis | Key-value store | CP/ACID |
| Riak | Key-value store | AP/BASE |
| SimpleDB | Wide-column store | AP/BASE |

# CRUD and REST

| CRUD(*)<br>(*) the four basic functions of persistent storage | HTTP(**)<br>(**) Methods for RESTful services | |
| --- | --- | --- |
| Create | POST | Create or add new entries |
| Read | **GET** | Read, retrieve, search, or view existing entries |
| Update | **POST** | Update or edit existing entries |
| Delete | DELETE | Delete/deactivate/remove existing entries |

# Map Reduce

# Scraping

amazon.com/dp/product-reviews/0316228532

Get the app

FREE 2-Day Shipping with Prime

**11 days left**

## amazon
Try Prime

All ▼

Departments ▼   Your Amazon.com   Today's Deals   Gift Cards & Registry   Sell   Help

Hello. Sign in
Account & Lists ▼   Orders   Try Prime ▼   Cart

The Casual Vacancy › Customer Reviews

## Customer Reviews

★★★☆☆ 5,801

3.2 out of 5 stars ▼

| | |
|---|---|
| 5 star | 26% |
| 4 star | 20% |
| 3 star | 17% |
| 2 star | 17% |
| 1 star | 20% |

### The Casual Vacancy
by J.K. Rowling

Format: Hardcover   Change

Price: $21.00 + Free shipping with Amazon Prime

Rate this item
☆☆☆☆☆   Write a review

Add to Cart

Add to Wish List

---

**Top positive review**
See all 2,678 positive reviews ›

1,395 people found this helpful
★★★★★ Brilliant, Disturbing, Not for Everyone
By Simply Keith on October 8, 2012

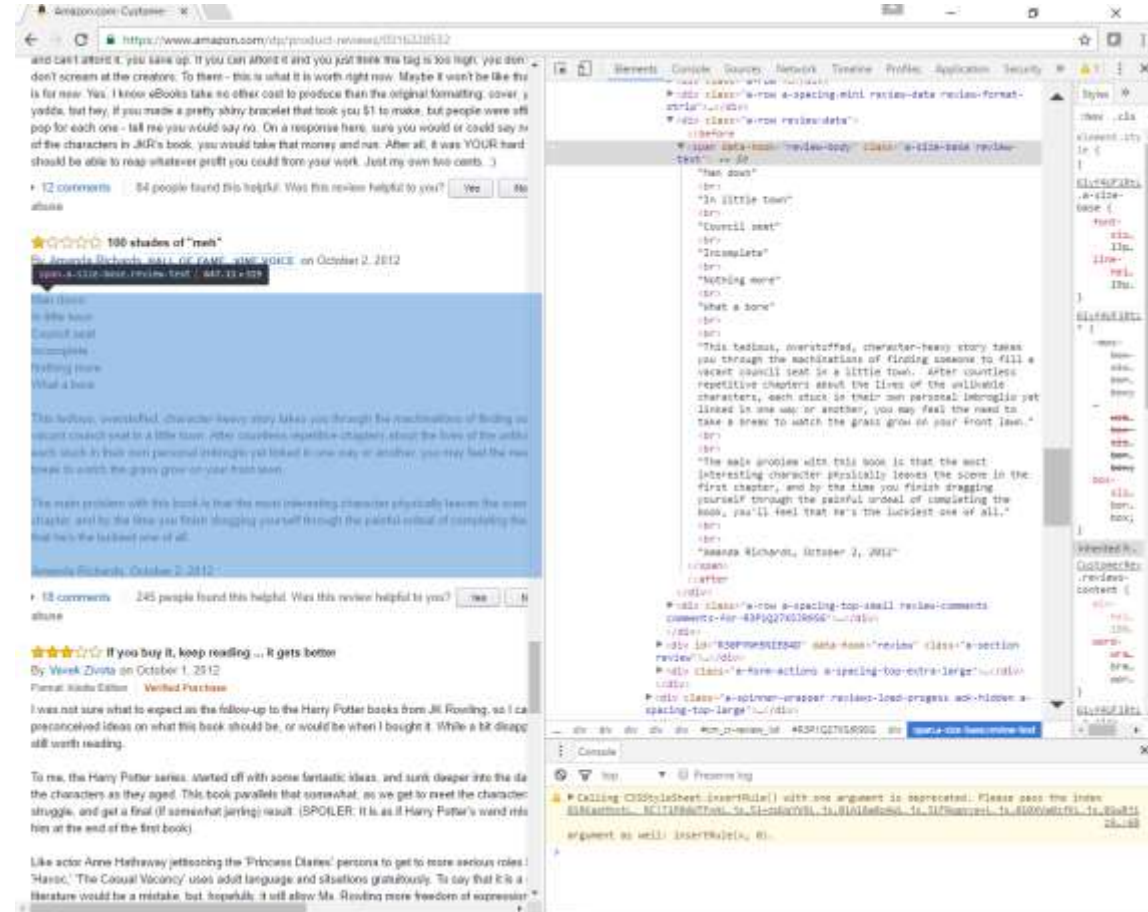Having read some of the more negative views, I have to say this: Would you have ranked "1984", "The Bluest Eye", "The Grapes of Wrath", or "Great Expectations" so badly? Guess what, some of the best stories aren't fun-filled light reading! Some of the best works are disturbing, even sad. So, if your view of literature is that a book can only be worthwhile if you can breeze through it having a fun adventure, then don't bother with this book. This is something else entirely.

That said, "The Casual Vacancy" is a disturbing character study. It is written in third-person omniscient point of view. It does require some effort to handle a story with some lift or so
Read more

**Top critical review**
See all 3,123 critical reviews ›

999 people found this helpful
★★★☆☆ From Someone Who's Actually Read It!
By Holly Day on September 27, 2012

Just to set the record straight on Casual Vacancy, I was given a copy as a gift. When you read the book, you need to forget this is the same author who gave us Mr Potter's world of magic. This is set in a tiny English town and deals with politics,class struggle,poverty,drug use,child abuse,rape,self-mutilation, suicide, pedophilia,mental illness and other ugly realities. Much to her credit the author does this with sharp comic wit, however it does all go on a bit too long with the ending being somewhat predictable and heavy-handed. It goes from being a lively comedy of manners to over-wrought slog by the end. I will say that I did enjoy her attempt to combine cutting comedy with social commentary. This is definitely not a book for the children in your house with the salty, profane language and delicate subject matters.
Read more

---

Sort by:  Top ▼      Filter by:  Verified pu... ▼   All stars ▼   All formats ▼      🔍 Keyword   Search

Showing 1-10 of 4,107 reviews (verified purchases). Show all reviews

★★★★★ Brilliant, Disturbing, Not for Everyone
By Simply Keith on October 8, 2012
Format: Hardcover   Verified Purchase

Having read some of the more negative views, I have to say this: Would you have ranked "1984", "The Bluest Eye", "The Grapes of Wrath", or "Great Expectations" so badly? Guess what, some of the best stories aren't fun-filled light reading! Some of the best works are disturbing, even sad. So, if your view of literature is that a book can only be worthwhile if you can breeze through it having a fun adventure, then don't bother with this book. This is something else entirely.

**Customers also viewed these items**

Fantastic Beasts and Where to Find Them: The Original Screenplay
by J.K. Rowling
$14.99
★★★★☆ 286

The Cuckoo's Calling (Cormoran Strike)
by Robert Galbraith

18

# Google Chrome's inspect element feature is a great tool to scrape web sites

# Application Programming Interfaces (APIs)

Slides © 2016 Ivan Corneillet Where Applicable
Do Not Reproduce Without Permission