

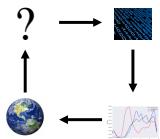
Review

DS 100, Spring 2017

Slides by:

Joe Hellerstein

hellerstein@berkeley.edu



Format of Exam

- When: Thursday 5/11 3PM
- Where: Genetics and Plant Bio 100
- How: Very similar to midterm in structure
- Cheat Sheets: 2!
- Topic areas: See syllabus.

Topics Covered in This Deck

- Relational Algebra and SQL
- Pandas
- Big Data: Data Warehousing, Map-Reduce, Spark
- Data Wrangling

Context and Connections

- We have seen a variety of data manipulation tasks
 - Both data "wrangling" and analysis
 - Diverse languages and frameworks
 - Pandas, SQL, MapReduce, Spark
- Check your fundamentals, study commonalities!
 - Matrix and set operations
 - Relational Algebra & Map/Reduce
 - Can you translate the operators from one into the other?
 - What operations appear frequently?
 - E.g. joins, groupby/sum/product, filtering on rows and columns
 - Relational algebra and linear algebra!

Relational Algebra & SQL

Relational Model & SQL

- Key Ideas:
 - Schemas, instances, queries, views, insert/delete/update
 - Relational algebra operators
 - SQL as a declarative language: everything referenced by value
 - "Not in" as in arrays or databases!
 - Many big data features under the covers: parallel execution, user-defined functions and aggregates
- What you should know:
 - Relational/SQL basics
 - Relational Algebra and basic equivalences
 - Basic SQL structure and syntax. Recognize incorrect logic (as opposed to syntax errors)
- What you don't need to know:
 - You will not be asked to compose queries from scratch
- How to study?
 - Review HW4, SQL notebook files from lecture

Relational Model

- Key Idea:
 - Schema: describes attribute (column) names and types, keys
 - Instance: a set (unordered, duplicate-free) of tuples (rows)
- Review the Boat Club from class

```
CREATE TABLE sailors(sailorId integer PRIMARY KEY, name text, rating integer, age float);
CREATE TABLE boats(boatId integer PRIMARY KEY, home text, color text);
CREATE TABLE reservations(reservationId integer, boatId integer, day date,
                         PRIMARY KEY (reservationId, boatId), day,
                         FOREIGN KEY (sailorId) REFERENCES sailors,
                         FOREIGN KEY (boatId) REFERENCES boats);
```

Note: primary keys underlined

Reserves

sailorId	boatId	day
28	101	10/19/12
58	103	11/12/12

Sailors

sailorId	name	rating	age
101	Intellect	blue	35.0
102	Intellect	red	35.0
104	Marine	red	35.0
103	Clipper	green	35.0

Boats

boatId	name	color	
31	Indigo	9	35.0
31	Jubilee	8	35.5
44	Guppy	5	35.0
58	Ruby	10	35.0

Relational Algebra Operators

Unary Operators: operate on **single relation instance**

- **Projection (π)**: Retains only desired columns (vertical)
- **Selection (σ)**: Selects a subset of rows (horizontal)
- **Renaming (ρ)**: Rename attributes and relations.

Binary Operators: operate on **pairs** of relation instances

- **Union (\cup)**: Tuples in $r1$ or in $r2$.
- **Intersection (\cap)**: Tuples in $r1$ and in $r2$.
- **Set-difference ($-$)**: Tuples in $r1$, but not in $r2$.
- **Cross-product (\times)**: Allows us to combine two relations.
- **Joins (\bowtie , \bowtie^*)**: Combine relations that satisfy predicates

As you review this material, do you see how it relates to SQL, Pandas and Spark?

SQL Data Model

- Basically relational, but "multisets"
 - I.e. can have duplicate rows in a table, and they matter
 - E.g. for COUNT, AVG, etc.
 - I.e. some tables have no primary key!
- FOREIGN KEYS

SQL

A declarative language: say "what" you want in the output, not "how" to get it.

DDL	DML
<code>CREATE TABLE <col> <type>, ...;</code>	<code>SELECT ...</code>
<i>Understand PRIMARY KEY, FOREIGN KEY clauses.</i>	<code>FROM ...</code>
	<code>[WHERE ...]</code>
<code>INSERT INTO TABLE values (<v>, <v>, ...);</code>	<code>[GROUP BY ...]</code>
<code>INSERT INTO TABLE SELECT ...;</code>	<code>[HAVING ...]</code>
	<code>[ORDER BY ...]</code>
<code>UPDATE TABLE SET ... WHERE ...;</code>	<code>[LIMIT ...];</code>
<i>note: value-based references in UPDATE!</i>	

SQL cheat sheet

Basic Queries

```
SELECT <col> FROM <table> WHERE <cond>
SELECT <col> FROM <table> WHERE <cond> AND <cond>
SELECT <col> FROM <table> WHERE <cond> OR <cond>
SELECT <col> FROM <table> WHERE <cond> NOT <cond>
SELECT <col> FROM <table> WHERE <cond> BETWEEN <val> AND <val>
SELECT <col> FROM <table> WHERE <cond> IN (<val>, <val>, ...)
```

Joins

The Joy of JOINs diagram illustrates four types of joins:

- LEFT OUTER JOIN**: All rows from table A, even if they do not exist in table B.
- INNER JOIN**: Only rows from both tables that have matching values in their join columns.
- RIGHT OUTER JOIN**: Only rows from table B, even if they do not exist in table A.
- CROSS JOIN**: Every row in table A paired with every row in table B.

Updates on JOINed Queries

You can use JOIN in your UPDATE:

```
UPDATE <table> SET <col> = <expr> WHERE <cond>
UPDATE <table> SET <col> = <expr> WHERE <cond> AND <cond>
UPDATE <table> SET <col> = <expr> WHERE <cond> OR <cond>
UPDATE <table> SET <col> = <expr> WHERE <cond> NOT <cond>
```

Useful Utility Functions

Common string functions:

- `TD_DATE` (Oracle, PostgreSQL, MySQL): Converts a timestamp to a date.
- `CONCAT` (MySQL): Concatenates multiple strings.
- `COALESCE` (PostgreSQL, MySQL): Returns the first non-null value.
- `CURRENT_TIMESTAMP` (MySQL): Returns the current timestamp.
- `MIN` / `MAX` (MySQL): Returns the minimum or maximum value in a column.
- `UNION` / `EXCEPT` (MySQL): Combines the results of two queries.
- `SELECT <col> FROM <table> WHERE <cond>` (MySQL): Returns data from both queries.
- `SELECT <col> FROM <table> WHERE <cond> OR <cond>` (MySQL): Returns data from the second query if the first one is not present.
- `REPORTING` (MySQL): Reports data from both queries.

Data Modification

Common update patterns with WHERE clause:

- `UPDATE <table> SET <col> = <expr> WHERE <cond>`
- `UPDATE <table> SET <col> = <expr> WHERE <cond> = <val>`
- `INSERT INTO <table> (ID, FIRST_NAME, LAST_NAME) VALUES (<val>, <val>, <val>)`
- `DELETE FROM <table> WHERE <cond>`
- `DELETE FROM <table> WHERE <cond> = <val>`
- `INSERT INTO <table> (FIRST_NAME, LAST_NAME) SELECT <col>, <col> FROM <table>`

Views

A **VIEW** is a virtual table, which is a result of a query. They are used to create virtual tables of complex queries.

```
CREATE VIEW <view_name> AS
  SELECT <col> FROM <table>
  WHERE <cond>
```

Indexes

Creating indexes on columns used in WHERE clauses:

```
CREATE INDEX <index_name> ON <table>(<col>)
```

Don't forget:

- Avoid overlapping indexes.
- Index on frequently updated columns.
- Indexes can speed up DELETE and UPDATE operations.

Reporting

Use aggregate functions:

- `COUNT`: counts the number of rows
- `SUM`: calculates the values
- `AVG`: calculate the average of the group
- `MIN / MAX`: smallest / largest value

Also

- **TABLESAMPLE BERNOUILLI(...)**
- Scalar functions (like "map") in SELECT and WHERE
 - Arithmetic, Math functions, string functions...
- Aggregate functions (like "reduce") in SELECT and HAVING
 - count(), sum(), average(), stddev()
 - Note: DISTINCT clause inside aggregates
- Table functions: generate_series()

Pandas and Numpy

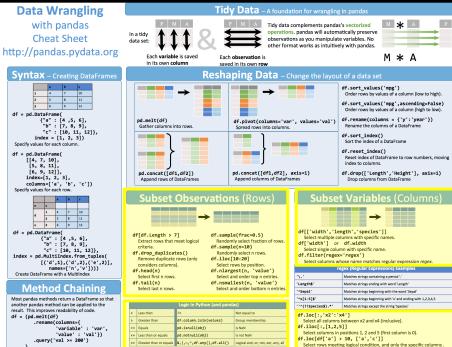
Arrays and DataFrames

- Key Ideas:
 - Array is a multi-d array; single type for all cells
 - Pandas DataFrame is like a relation (table); single type per column
 - Unlike relational model, fixed order to rows (allowing array-style indexing)
- What you should know:
 - Basic numpy array handling
 - Key pandas operations (see next slides)
 - Recognize incorrect logic
- What you don't need to know:
 - You will not be asked to write code from scratch
- How to study?
 - Review HW1-3, notebook files from lecture
 - Be able to translate basic SQL block above into Pandas

Topics

- Data Frame & Numpy Arrays
- Subsets
- Groupby/Agg
- Joining

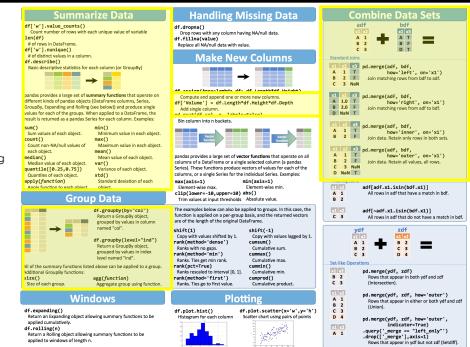
https://github.com/pandas-dev/pandas/blob/master/doc/cheatsheet/Pandas_Cheat_Sheet.pdf



Topics

- Data Frame & Numpy Arrays
- Subsets
- Indexing
- Groupby/Agg
- Joining

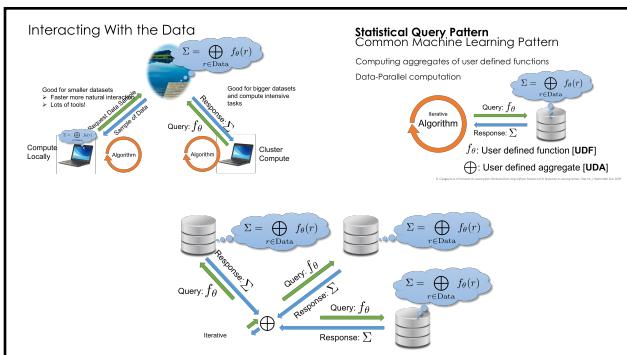
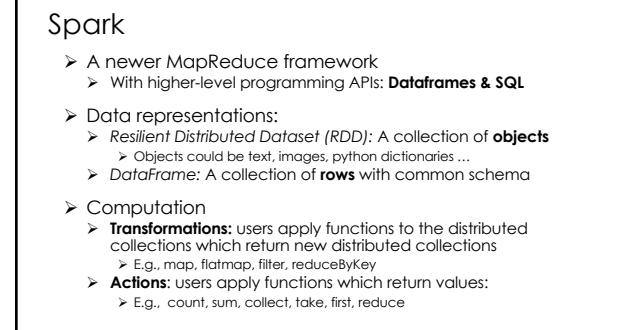
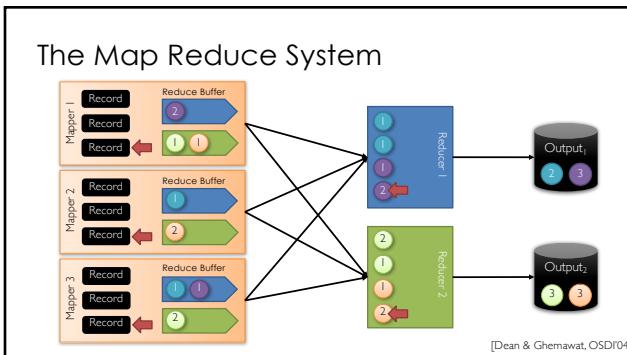
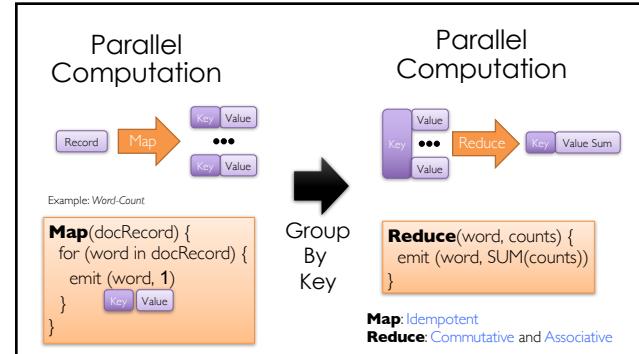
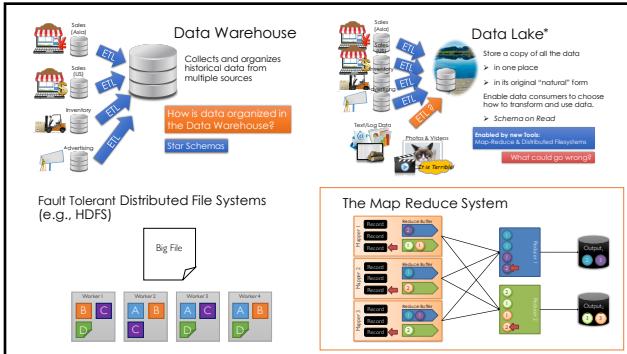
https://github.com/pandas-dev/pandas/blob/master/doc/cheatsheet/Pandas_Cheat_Sheet.pdf



Big Data

Big Data

- Key Ideas:
 - Schema on load (data warehousing) vs. on use (data lakes)
 - Fault-tolerant file systems
 - Data-parallel computing (e.g. MapReduce)
 - Statistical Query Pattern
- What you should know:
 - Star schema design
 - File system data layouts and failure handling
 - Map vs. Reduce
 - Spark Transformations vs. Actions
 - How to compute basic linear algebra operations in statistical query pattern
- What you don't need to know:
 - You will not be asked to write code from scratch
 - Fault tolerance in MapReduce
- How to study?
 - Review HW 7 and lecture notebooks



Data Wrangling

- Key Ideas:
 - Wrangling output determines your downstream analytics
 - Assessing data quality like querying + vis
 - Data transformation also like querying!
 - A petit-dish for many of the topics in class
- What you should know:
 - Structural issues: arrays, relations, and converting between Granularity; assessing it (keys) and changing it (grouping)
 - Faithfulness: finding and handling outliers
 - Basic UNIX commands
- What you don't need to know
 - UNIX command flags
 - Trifecta Wrangler
 - Memorizing "Structure, Granularity, Faithfulness, Temporality, Scope" (that's just to organize the ideas)
 - Live-coding Pandas/Python details like "read_csv" or strip() or head()
- How to study?
 - Lecture notes
 - HW3

Rough Guide to Topics

- **Structure:** the "shape" of a data file
- **Granularity:** how fine/coarse is each datum
- **Faithfulness:** how well does the data capture "reality"
- Temporality: how is the data situated in time
- Scope: how (in)complete is the data

We didn't really discuss temporality/scope much

Structure

- Relations and Relational languages (algebra, SQL)
- Arrays and Linear Algebra
- Coarse structure
 - Record and field delimiters
 - Common schemas
 - Nested?
- Value encoding: nominal vs. ordinal vs. quantitative
- See connections to SQL and DataFrames, n-d Arrays

Granularity

- Keys as an indicator of granularity
- Granularity can be coarsened via groupby
 - Grouping columns become a primary key of result
- Granularity can be refined via join
 - With more detailed tables
 - Foreign keys
 - Understand the # of matches for various joins

Faithfulness

- Data in context
 - Application context
 - E.g. human beings are unlikely to be 1000 years old
 - Constraints: relation to "types" or "domains"
 - E.g. "123456789" is not a valid zip code
 - E.g. the word "spunow" is not in the english dictionary
 - Cross-record context
 - Detecting Outliers
 - "Center": e.g. average and median
 - "Spread": e.g. stdev, inter-quartile range
 - Resolving Outliers
 - Trimming: setting outlier values to N/A or NULL
 - Winsorizing: setting outlier values to the nearest non-outlier value

