# Analytics-R'-Us

## Schema Integration and Justification Team

Demand Prediction Analysis

DSE 203 Presentation #2

10/28/2017

**Team:**

Josh Wilson

Amisha Bhanage

Ken Kroel

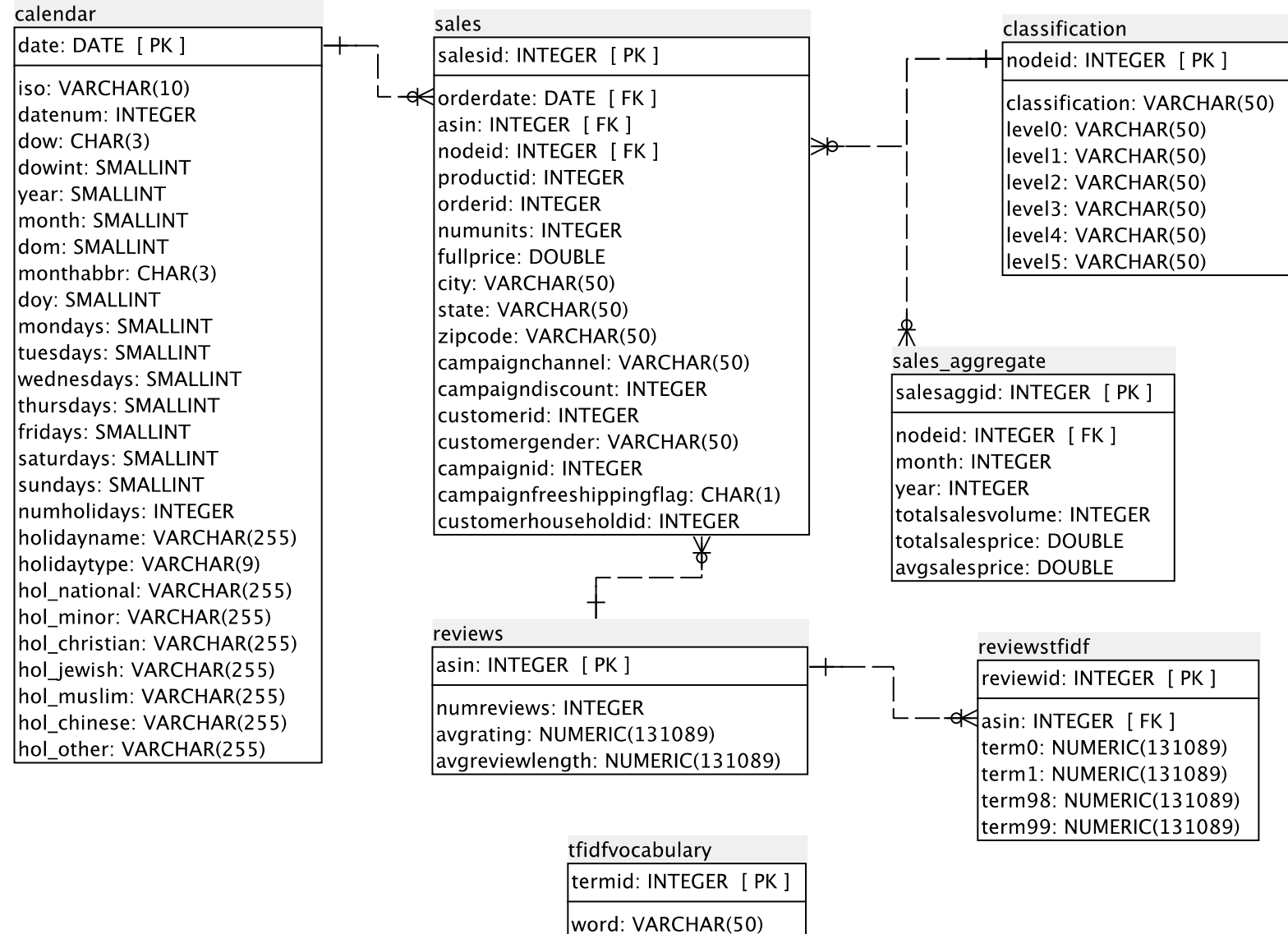Mai Huynh

# Integrated Schema Provides

- Sales Information

| | |
|---|---|
| Product ID (Books) | Unit Price |
| Product Classification (flattened) | Average Price |
| Order ID | Number of Units |
| Customer ID | Campaign ID |
| Customer Gender | Campaign Discounts |
| Period (Day, Week, Month, Year, Holidays) | Campaign Channels |
| Geography (State, City, Zip) | Free Shipping Flag |

- Book Review Information

| | |
|---|---|
| Review Text TF-IDF of 100 word | Number of Reviews |
| Average Rating | Average Review Length |

# Virtualized Integrated Schema

# Virtualized Integrated Schema Mapping

- Sales → PostgreSQL

- Calendar → PostgreSQL

- Classification → flattened info from JSON / AsterixDB

- Reviews → ASIN/product level summary from JSON / AsterixDB

- ReviewsTfIdf → TF-IDF info for top 100 vocabulary terms for each review

- TfIdfVocabulary → maps each vocabulary word to its numeric term

- Sales_aggregate → summary of sales volume and price by classification, month, and year

- Full mappings between integrated schema and source data available on Github

# Example of Query Decomposition

- Query:
  - Total sales volume for "Computer Programming" in in July-2013
- Source data:
  - Sales ( orderdate, nodeid, numunits, orderdate) ← PostgreSQL
  - Classification (nodeid, classification) ← JSON from AsterixDB
- Query breakdown

**SELECT** agg.nodeid, agg.MONTH, YEAR, agg.total_sales_volume

**FROM** sales_by_month agg **INNER JOIN** classification cls

**ON**      agg.nodeid      = cls.nodid

**WHERE**  MONTH      = 7 **AND** YEAR      = 2013

**AND** classification = "for computer programming"

# Additional Query Ideas for Stakeholders

1. What is the average rating for each book?

2. What is the average rating of books within each category?

3. What is the variance in ratings within each category?

4. What is the geographic distribution of orders for each category?

5. What is the variance in number of orders by book in each category? (i.e. do some books within a category sell more than others?)

6. What is the variance in number of orders in each category by month/quarter? (i.e. is demand seasonal?)

7. What percentage of sales in each category are associated with a campaign? With each specific type of campaign?

8. How correlated is demand with number of reviews, average review rating, discount percentage, days until holiday season?

9. Are specific customers more influenced by certain types of campaigns?

10. Are products from the same category more likely to be bought in the same order?

# Q & A

Thank you!