# Project 1 Checkpoint

**Sample Superstore Dataset**

# Columns I'll be working on:

- **Measures:** Sales, Profit, Quantity, Discount
- **Dimensions:** Region, State, Segment, Sub-Category, Category, Ship Mode
- **Derived fields I'll create:**
  - profit_margin = Profit / Sales (guard Sales=0)
  - is_loss = Profit < 0
  - big_discount = Discount >= 0.20
  - discount_tier = categorized as "None", "Low", "Medium", or "High"

# Calculations:

## Calculation 1: Average profit margin by sub-category within each region

- **Uses:** Profit, Sales, Sub-Category, Region
- **Method:** For each (Region, Sub-Category), compute weighted margin = sum(Profit) / sum(Sales)
- **Output file:** margin_by_region_subcategory.csv
- **Columns:** Region, SubCategory, total_sales, total_profit, profit_margin

## Calculation 2: Loss rate for high-discount lines by state and segment

- **Uses:** Discount, Profit, State, Segment
- **Method:** Filter rows where Discount >= 0.20; for each (State, Segment), compute percent with Profit < 0
- **Output file:** loss_pct_high_discount_by_state_segment.csv
- **Columns:** State, Segment, num_lines, num_losses, loss_pct

## Calculation 3: Regional Performance by Customer Segment

- **Uses:** Region, Segment, Sales, Quantity
- **Method:** For each (Region, Segment) combination, calculate average order value = total_sales / total_quantity
- **Output file:** avg_order_value_by_region_segment.csv
- **Columns:** Region, Segment, total_sales, total_quantity, avg_order_value

## Calculation 4: Discount Impact on Order Size by Category

- **Uses:** Discount, Quantity, Category, Sales
- **Method:** Create discount tiers (None: 0%, Low: 0-20%, Medium: 20-40%, High: 40%+); for each (discount_tier, Category), calculate average quantity and average sales
- **Output file:** discount_impact_by_category.csv
- **Columns:** discount_tier, Category, num_orders, avg_quantity, avg_sales

# Team:

David Vargas & Alberto Puliga

# Diagram: