

# Our datasets

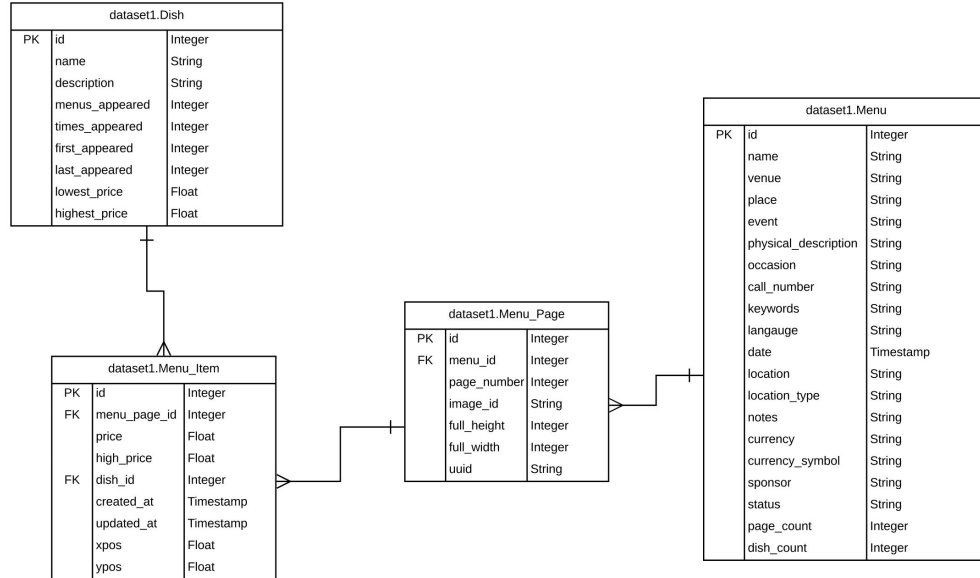
## **Dataset 1 - New York Public Library Restaurant Menus**

- Tables:
  - Dish
  - Menu Item
  - Menu Page
  - Menu

## **Dataset 2 - Yelp Review Data**

- Tables:
  - Business
  - User
  - Review

# ERD - v1



# Our questions

- Were ratings of a restaurant the same between a user and their friends?
- Does having more items on a menu correlate with lower ratings?
- Does having fewer dishes on a menu page correlate with higher ratings?
- Do ratings of a restaurant have anything to do with whether the menu is vertical or horizontal?

# Dataset1

## Transformations

### Apache Beam:

- Price range of food item
  - Calculate range between price and high price for a menu item
- Orientation of Menu
  - Determine whether a menu is horizontally or vertically oriented
- Average Number of Dishes per Page of Menu
  - Calculate the average number of dishes per page from the Menu's dish count and page count

### SQL:

- Currency table
  - Remove currency information from Menu, create new table
- Location table
  - Remove location information from Menu, create new table

# Dataset2 Transformations

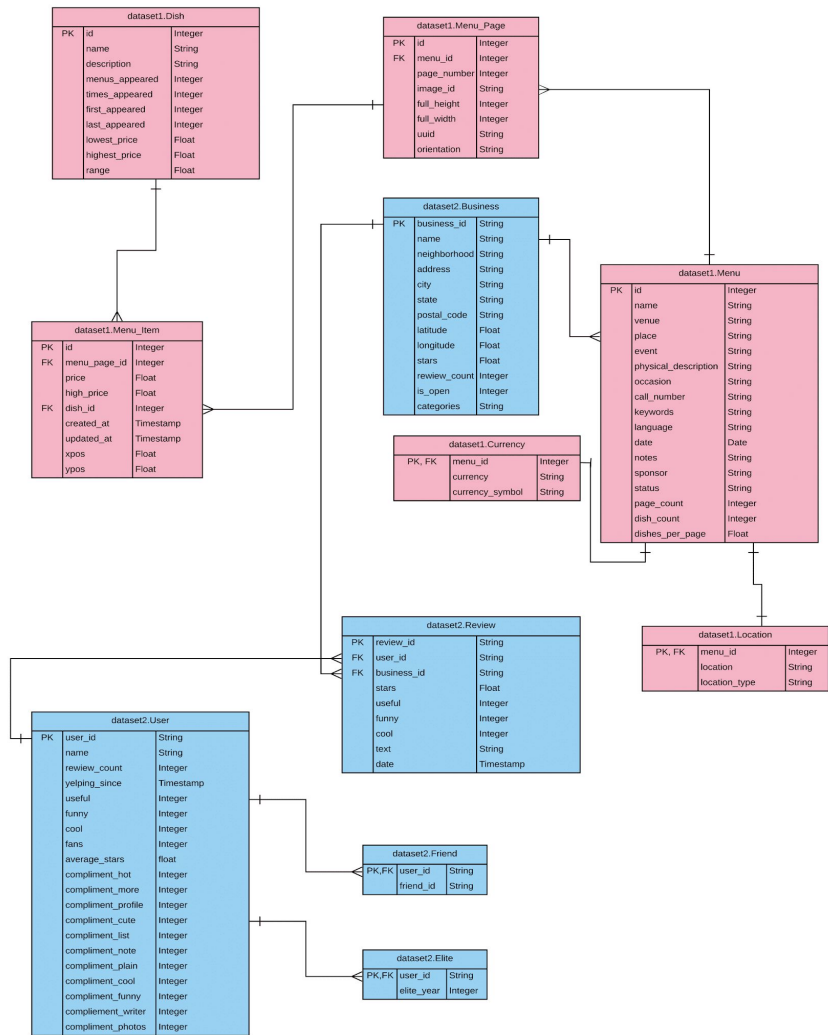
## Apache Beam:

- Business Table
  - Transform Business name to match the formatting in the NYPL dataset
- Elite Table
  - Transform Elite Year list in User into its' own table
- Friend Table (demo)
  - Transform Friend list in User into its' own table
  - <https://bit.ly/2V8k8Ap>

## SQL:

- User table changes
  - Remove Friend & Elite

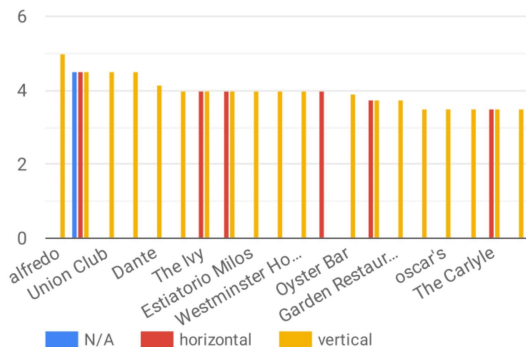
# ERD - v5



# Cross-Dataset Findings

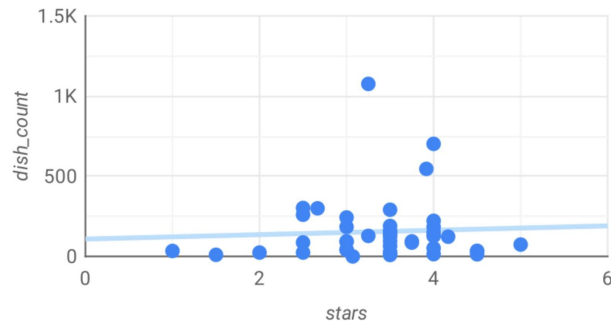
```
SELECT b.newname, b.stars, mp.orientation
FROM `dogwood-theorem-230119.dataset1.Menu` m
JOIN `dogwood-theorem-230119.dataset2.Business_1` b ON UPPER(m.place) = UPPER(b.newname)
JOIN `dogwood-theorem-230119.dataset1.Menu_Page` mp ON mp.menu_id = m.id
ORDER BY b.stars DESC
```

Menu Page Orientation vs. Average Star Rating for the Restaurant



```
SELECT b.newname, m.dish_count, b.stars
FROM `dogwood-theorem-230119.dataset1.Menu` m
JOIN `dogwood-theorem-230119.dataset2.Business_1` b ON UPPER(m.place) = UPPER(b.newname)
WHERE dish_count is not null
```

Dishes on a Menu vs. Average Star Rating for the Restaurant



# Airflow

```
# DAG section
with models.DAG('workflow',
                schedule_interval=datetime.timedelta(days=1),
                default_args=default_dag_args) as dag:

    # Beam Tasks
    business_beam = BashOperator(
        task_id='business_beam',
        bash_command='python /home/stoddartisabel/' + business_script)

    elite_beam = BashOperator(
        task_id='elite_beam',
        bash_command='python /home/stoddartisabel/' + elite_script)

    friend_beam = BashOperator(
        task_id='friend_beam',
        bash_command='python /home/stoddartisabel/' + friend_script)

    # SQL Task
    create_user_table = BashOperator(
        task_id='create_user_table',
        bash_command=sql_cmd_start + ''' + sql_user + ''')

    [business_beam, elite_beam, friend_beam] >> create_user_table
```



# Future improvements

This solution might be improved by:

- More comprehensive location data for our restaurants, looking at demographic information for neighborhoods
- More information for each Menu from NYPL
- More current menus
  - A lot of our data was from old menus, which wouldn't exist in the Yelp database since Yelp was founded in 2004