# Lecture 8:
# DBSCAN Clustering and Analytics Dataflow

## COSC 526:
## Introduction to Data Mining
## Spring 2021

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE
KNOXVILLE

**BIG ORANGE. BIG IDEAS.**®

# Lecture Outline

- Assignment 8
  - Three problems with missing data
- Project Discussion
  - Define key questions and tentative title
- Dataflow and DBSCAN
  - More in the next slides
- If time left
  - Live chat and video - https://www.youtube.com/watch?v=_2u_eHHzRto

BIG**ORANGE**
BIG**IDEAS**

# Lecture Outline

- Use work in paper "*Development of a Scalable Method for Creating Food Groups Using the NHANES Dataset and MapReduce*" to learn about:
  - Using a dataflow for tackling a data problem
  - Structuring a research project into a set of slides
    - Motivation, goals, background, methodology and results
  - Using a different clustering method than k-mean
    - How to cluster numerical data with the Density-based spatial clustering of applications with noise (DBSCAN)
    - How to set up the setting parameters of the DBSCAN
  - Using code from other scientists
    - Re-use rather than rewriting from scratch
    - Replicability of work as the first step for new research

# Assignment 8

# Relevant Open-source Dataset

- Use a dataset with well-known and broadly used data format

    **NHANES:** National Health and Nutrition Examination Survey

    - Medical, demographic, and dietary records
    - Available to the public for free
    - *Contains subjective food groups provided by USDA*

**BIG ORANGE**
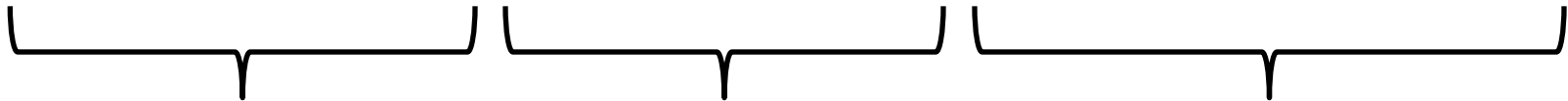**BIG IDEAS**

# NHANES Dietary Data

- Dietary intake of 64,653 Americans

- 7,494 unique food items

- 1,587,750 food entries

- 46 nutrient features for each food item
  - Macronutrients (e.g., fats, carbohydrates)
  - Micronutrients (e.g., vitamins, minerals)

# NHANES Dietary Data

- Dietary intake of 64,653 Americans
- **7,462 unique food items**
- 1,587,750 food entries
- 46 nutrient features for each food item
  - **Macronutrients (e.g., fats, carbohydrates, proteins)**
  - Micronutrients (e.g., vitamins, minerals)

BIG **ORANGE**
BIG **IDEAS**
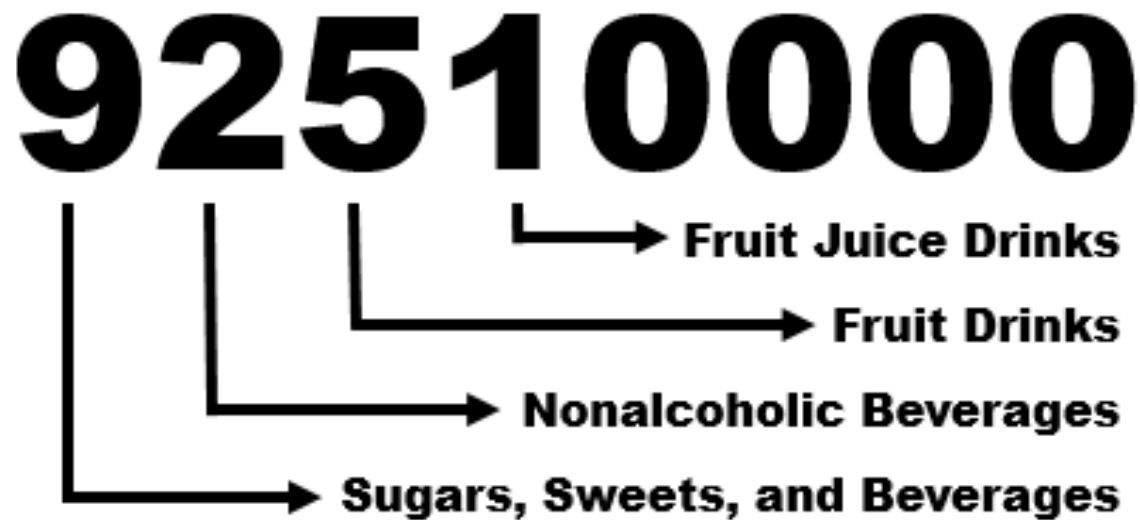
# Structure of Dietary Data Item

<143672, 92510000, 3, 0, 8:15am, 7, 3, 10.1, 4, 3.45, 10, 178, …>

- Participant ID
- **USDA Food Code**

- Meta Data

- **Macronutrients**
- Micronutrients

# USDA Food Classification

- Subjective and general
- Categorical, not nutrient-driven

## 92510000

- → **Fruit Juice Drinks**
- → **Fruit Drinks**
- → **Nonalcoholic Beverages**
- → **Sugars, Sweets, and Beverages**

# Assignment 8 – Problem 1

- Data file: ./data/data-1.csv
  - Contain no missing values
- Problem: Cluster food items in the file based on carbohydrate and fat:
  - *Use the k-Means from the Spark MLlib to cluster data points*
  - *Determine the optimal value for k using the* **elbow method**
- Metric of quality:  Use <u>Within Set Sum of Squared Errors</u>
  - This method is built into the Spark kMeans model and can be accessed with model.computeCost()
- Note: the clusters provide a ground truth for comparison to clusters we find later using data with missing values.

BIG**ORANGE**
BIG**IDEAS**

# Assignment 8 – Problem 1

Steps:

- Define the optimal value for K

- Cluster food items (using K)

- Plot clusters by *fat* and *carbohydrate* content

# Assignment 8 – Problem 2

- Data file: ./data/data-2.csv
  - Missing values for the carbohydrate content of some food items
  - The data were removed from a specific set of food items (i.e., food items with carbohydrate value near 0.5)
- Define a method to remove food items with missing any  macronutrient values and apply this method to the data.
- Cluster the modified data and plot the results
  - Use the same K value as in Problem 1
- Note: we provide you with the code that loads the data and reports the percentage of values missing for each macronutrient (i.e., carbohydrates and fat)

# Assignment 8 – Problem 3

- Data file: ./data/data-3.csv

  - Missing values for the fat or carbohydrate content of some food items (but not both for a single food item).

  - The data were removed from a food items randomly

PART 1:

- Define and apply a method to fill missing values with the mean of other values

  - *E.g., for missing values in fat, fill with the mean of fat values that are present*

- *Cluster the modified data and plot the results*

  - Use the same K value as in Problem 1

# Assignment 8 – Problem 3

*PART 2:*

- *Use the code for Problem 2 to remove data with missing values rather than filling the gaps (as you did in PART 1)*

- *Cluster the modified data and plot the results*
  - Use the same K value as in Problem 1

BIG**ORANGE**
BIG**IDEAS**

# Assignment 8 – Problem 4

- **Observe and describe:** Can you summarize your findings in each problem? Can you compare and contrast the findings across problems? How did each method for dealing with missing data (i.e., remove of filling) change the clustering outcome?

- **Impact of K:** What value did you choose for K in Problems 1-3? You based the selection of your K on the first dataset (i.e., no missing data). Do you expect a different value of K if you had used the elbow method with the second or third dataset? If yes, propose changes to your current solutions.

# Assignment 8 – Problem 4

- **Building assumptions on data distributions:** Now look at the plot of clusters in Problem 1. Logically, there cannot be more than 1 gram of (carbohydrate + fat) in 1 gram of food. In your plot this can be seen in the form of a diagonal line from the top-left to bottom-right (where the sum of fat and carbohydrate content is equal to 1). How can you use this information to improve the way you fill missing values? Can you think of other methods to fill missing values? (HINT: logistic regression)

# Project Steps

# Project (I)

- Step 0: search for datasets
  - Discussion in class of datasets identified

# Assignment 6 – Problem 5 DONE!

**Find an interesting dataset.**

- Is the dataset available? Do you have the dataset source (with url)?
- It your dataset large enough to allow for interesting analysis and non-trivial results?
- How were the data were collected?
- What is the significance of the data?
- Describe the number of rows, objects, or data points
- What information is contained in each rows, objects, or data points?
- What is the data types (int, str, char, float, etc.) and numerical ranges where appropriate
- What file(s) do you need to parse (if multiple files are available)?

# Project (II)

- Step 1: Address (and answer) these questions
  - Do you work alone or in team? If in team, indicate who is your collaborator and what are the skills / expertise he/she is bringing to the project
  - What is the dataset you are considering?
  - What are the possible key question(s) you want to answer? Are the questions too general? Are the questions too narrow?

# Assignment 7 – Project Questions DONE

**Answer the following questions, in a couple sentences each, in the cells provided below. Note that you can indicate that you do not have a dataset or/and a set of questions, and you wish to be helped in completing this task.**

- Do you wish to work alone or in team? If in team, indicate what are the skills / expertise he/she/they can bring to the project. Suggestion: you may consider to work in team with another student who works in a different scientific domain.

- Have you identified your dataset? If yes, what is the dataset you are considering?

- Have you identified possible key question(s) you want to answer? If yes, list the questions.

# Project (III)

- What are the scientific questions you are answering with your tool/framework? Be as specific as possible.

- What is the tentative title of your project?

- What are the milestones you want to meet from now until the end of the semester when you will present your poster at the poster showcase? Be as specific as possible. Write the date of the deadline and the task(s) you want to achieve for that deadline.

BIG**ORANGE**
BIG**IDEAS**

# Assignment 8 – Project Questions <span style="color:red">THIS WEEK</span>

**Answer the following questions, in a couple sentences each, in the cells provided below.**

- What are the scientific questions you are answering with your tool/framework? Be as specific as possible.

- What is the tentative title of your project?

- What are the milestones you want to meet from now until the end of the semester when you will present your poster at the poster showcase? Be as specific as possible. Write the date of the deadline and the task(s) you want to achieve for that deadline.

# Project (IV) NEXT WEEK in Assignment 9

**Motivation** Describe the motivation of your work. To build the motivation, you can answer these questions:

- What is the problem you are tackling?
- How is the problem solved today?
- Write a paragraph of 200 - 300 words

**Contributions** List between 2 and 4 contributions of your work. Contributions are bullet points that define your solution. E.g.,

- We build a system that ....
- We validate the system accuracy by ....
- We measure the performance of the system by ...
- Write a section of 150 - 200 words

# Project (IV) NEXT WEEK in Assignment 9

**Tests** List the type of tests (measurements) you will perform. E.g.,
- What are your metrics of success?
- Where do you run your tests?
- What tests do you perform?
- How many times do you run each test?
- What do you measure?
- Write a section of 250 - 350 words.