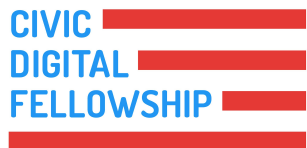# Using Machine Learning to Predict NAPCS Codes from Product Descriptions

**Economic Statistical Methods Division**
Carma Hogue —ESMD Division Chief
Andrea Roberson —ESMD Mathematical Statistician

CIVIC DIGITAL FELLOWSHIP

United States Census Bureau

**Beverly Bendix** & **Fatima Koli**
Carnegie Mellon University
M.S Public Policy & Management

Columbia University
M.S Data Science

# The project was motivated by item 22 of the Economic Census

## ITEM 22: DETAIL OF SALES, SHIPMENTS, RECEIPTS, OR REVENUE

Of the **$,000.00** of Sales, Shipments, Receipts, or Revenue reported in **Item 5**, what was the value for each product or service?

| Description | Value | Product Code |
|---|---|---|
| **1.** Retail sales of fresh meat and poultry<br>(Report deli meats on line 6a and meats sold in a frozen state on line 7a.) | $ ,000.00 | 5000025000 |
| **2.** Retail sales of fresh fish and seafood<br>(Report fish and seafood sold in a frozen state on line 7a.) | $ ,000.00 | 5000050000 |
| **3.** Retail sales of fresh fruit and vegetables<br>(Report frozen fruits and vegetables on line 7a.) | $ ,000.00 | 5000075000 |
| **4.** Retail sales of eggs and dairy (except ice cream)<br>(Report deli cheeses on line 6a and ice cream and other frozen dairy products on line 7b.) | $ ,000.00 | 5000100000 |

# 3.7 Million Businesses are Impacted

## ...equal to 6 tractor trailers
## of notification letters

CIVIC
DIGITAL
FELLOWSHIP

# The 1,200+ NAPCS codes are complex & can be ambiguous

## Which category?

- Retail sales of candy, prepackaged cookies, and **snack foods**

- Retail sales of food **dry goods** and other foods purchased for future consumption (Include flour, sugar, fats and oils, coffee, honey, jams and jellies, pasta, and crackers.)
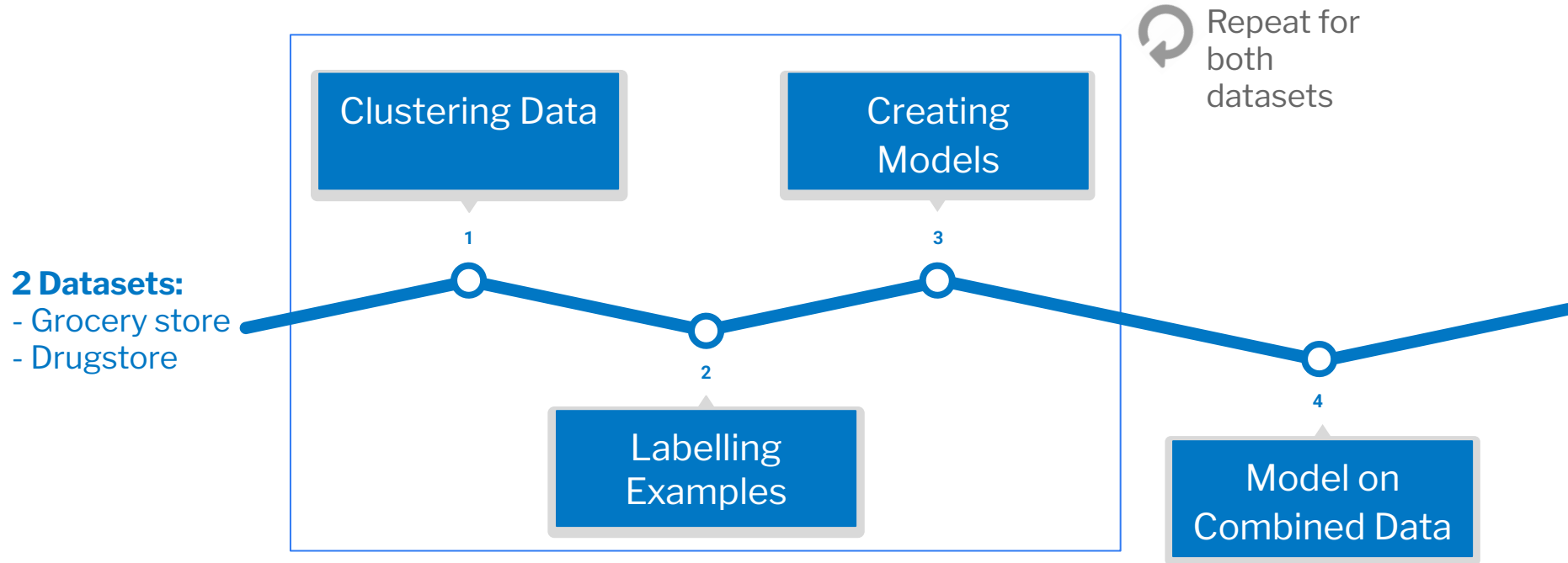
- Retail sales of fresh fruit and vegetables



Photo credit: Jeff Meisel

# This is a text classification problem

Using the UPC product descriptions, can we predict the NAPCS code?

# Process



Repeat for both datasets

Clustering Data

Creating Models

**2 Datasets:**
- Grocery store
- Drugstore

1

2

3

4

Labelling Examples

Model on Combined Data

CIVIC
DIGITAL
FELLOWSHIP

# Clustering Uncovers Patterns in the Data

| GREETING CARDS | HAIRCARE | CHOCOLATE |
|---|---|---|
| **Cluster2** | **Cluster 6** | **Cluster 22** |
| • ['card'] | • ['sham'] | • ['choc'] |
| • ['mark'] | • ['oz'] | • ['oz'] |
| • ['gift'] | • ['pantene'] | • ['milk'] |
| • ['vgc'] | • ['dry'] | • ['bar'] |
| • ['hmsignture'] | • ['shld'] | • ['dk'] |
| • ['sympthy'] | • ['head'] | • ['dark'] |
| • ['evryday'] | • ['color'] | • ['chip'] |

# ... But Clustering Results Can be Misleading

Descriptions closest to "MEOW MIX ORIGINAL CAT FOOD 3.15LB"

| Description | Similarity |
| --- | --- |
| FRISKIES PARTY MIX ORIGINAL  20OZ | .976 |
| CHEEZ-IT SNACK MIX ORIGINAL 4.5Z | .958 |
| CORN NUTS ORIGINAL  4OZ | .948 |
| CHIPS AHOY! ORIGINAL COOKIES 13OZ | .944 |

CIVIC
DIGITAL
FELLOWSHIP

# UPC Descriptions are Used to Predict NAPCS Codes

**UPC Description** $\longrightarrow$ **NAPCS Code**

SALMON KING HEAD WILD FRESH

Fresh fish and seafood

NO NONSENSE XSPORT NO SHOW
SOCK ASTD

Women's socks or Men's socks?

MINI TUB OF FUN GENERIC

?

CIVIC
DIGITAL
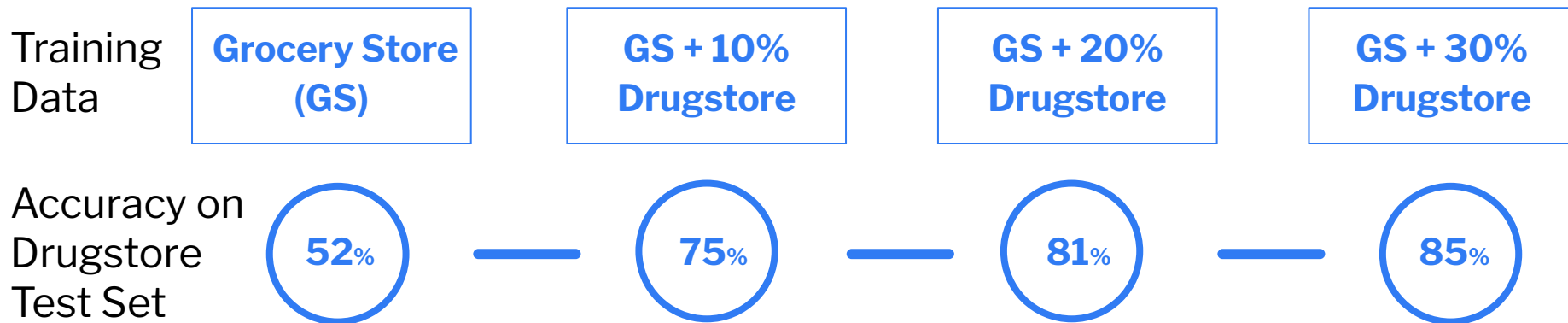FELLOWSHIP

# Label, Label, Label!

| Dataset: | Drugstore | Grocery Store |
|---|---|---|
| # Obs. : | 14K | 597K |

**3 difficult weeks later:**

| | Drugstore | Grocery Store |
|---|---|---|
| Labeled: | 14K | 36.7K |
| % of Dataset | 100% | 6.1% |
| NAPCS Codes: | 44 | 45 |

# Accuracy of Our Models:
**95%** on drugstore data
**90%** on grocery store data

# Tested Model's Ability to Generalize

Training Data

| Grocery Store (GS) | GS + 10% Drugstore | GS + 20% Drugstore | GS + 30% Drugstore |
|---|---|---|---|

Accuracy on Drugstore Test Set

52% — 75% — 81% — 85%

# This will significantly modernize the Economic Census

**BEFORE**

**High** respondent burden

Poor data quality

Collected **every 5 years**

**AFTER**

**Low** respondent burden

Automated classification is more **consistent** and **accurate**

Could be done **annually!**

CIVIC DIGITAL FELLOWSHIP

# Recognition

Thank you to **Andrea Roberson** for her machine learning acumen and guidance and to **Carma Hogue** for taking us onto her team!

CIVIC
DIGITAL
FELLOWSHIP