# Machine Learning
## Optimizing the Commodity Flow Survey

**Andrew Cukierwar**

Washington University in St. Louis

Mentor: James Hinckley, Economic Reimbursable Surveys Division

# Commodity Flow Survey

- Conducted **every 5 years**

- Co-sponsored by Census and BTS

- Each quarter, select businesses report data on up to 80 of their shipments

# SCTG Codes

Standard Classification of Transported Goods

- SCTG code for **each shipment**

- Finding code is **time-intensive**
  - Over 500 SCTG codes

- **Incorrect** SCTG codes

- Ex: 01003 – Live Poultry



One of 15 Pages in Booklet

# Solution
## Machine Learning

- Use **machine learning** techniques to **automate** assignment of SCTG codes to shipments
  - Classify shipments into SCTG codes
  - Allows survey responders to **save great amount of time**

- Implemented in Python using the Scikit-Learn library

# Classification

- Use X-variables to identify what class (Y-variable) an item belongs to

| | X1 | X2 | X3 | Y |
|---|---|---|---|---|
| **Item 1** | 5 | 1 | 3 | O |
| **Item 2** | 4 | 5 | 4 | + |
| **Item 3** | 2 | 7 | 6 | + |
| **Item 4** | 2 | 9 | 8 | O |
| **Item 5** | 3 | 7 | 6 | + |

# Initial Data Preparation

- Initially **3 million** rows of data

- Only **1.5 million** rows pass edits
  - <u>Ex:</u> Shipment is missing SCTG code

- Remove Duplicate Rows
  - Left with **230,000 unique** rows

### Rows of Data

# Text Preprocessing

- Data includes short **text description of shipment**

- Manipulate text to make it **uniform**

28 STEEEL BEAM,S
28 STEEEL BEAMS
STEEEL BEAMS
steeel beams
steel beams

# Text to Matrix

- **Must convert text to matrix form**
  - Algorithms only work on **numerical data**

- **TF-IDF** (Term Frequency – Inverse Document Frequency)
  - Bag of Bigrams

# Bag of Bigrams

Each description is split into a 'bag' of individual words (**unigrams**) and pairs of consecutive words (**bigrams**)

# TF-IDF

- Create matrix of counts of unigrams and bigrams

| | Description | steel | beams | pipe | steel beams | steel pipe | steel steel |
|---|---|---|---|---|---|---|---|
| 1 | steel beams | 1 | 1 | 0 | 1 | 0 | 0 |
| 2 | steel pipe | 1 | 0 | 1 | 0 | 1 | 0 |
| 3 | steel steel | 2 | 0 | 0 | 0 | 0 | 1 |

- Multiply matrix values by inverse of frequency in full body of text

- Resulting matrix is **large but sparse**
  - 230,000 Rows and 170,000 Cols = **39.1 billion values**
  - Nearly all values are **0's**

# Matrix Processing

- Add two columns of data to the matrix
  - Price per weight
  - Average word length

- **Normalize** all data to be between 0 and 1

# Logistic Regression Classifier

- **Fast** for large and sparse data
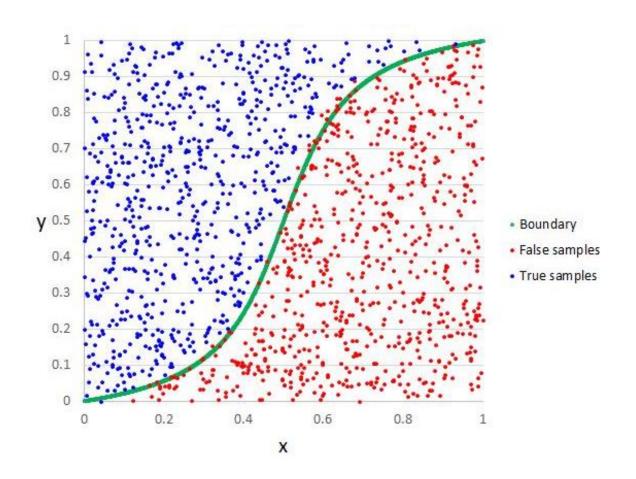
- **Highest accuracy** for this data

- Returns **probability** scores
  - Probability of X shipment having SCTG code Y is Z%

# Training and Testing

- Split data into training and testing data (90-10 split)

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Train | Train | Train | Train | Train | Train | Train | Train | Train | Test |

- Use training data to "create" model

- Use testing data to determine accuracy of model

# Model Process

# Results

- Predict **first 2 digits** of SCTG code with **65% accuracy**

- Predict **full SCTG code** with **48% accuracy**

- **Given first 2 digits** of SCTG code, can predict **full SCTG code** with **70% accuracy**

# Limitations

- Model is only as good as the data used
  - **Using survey respondent data to train**
  - Short descriptions of shipments (100 char limit)
  - Many classes
  - 43XXX – Mixed Freight
  - Codes for miscellaneous products

# Future Improvements

- ## Improve quality of data

  - Train model only using label-verified shipments

  - Ensure spelling of descriptions

  - Increase shipment description length

  - Train model with more data

- ## Use more advanced algorithms (deep learning)