# Putting the Cure into Strategic Procurement

NLP soft weighted voting ensemble multi-class classifier to predict spend categories using Government of California's 2012 -2015 purchase order data

Dene Stalk

BrainStation Data Science Diploma Capstone Project

March 2020

# Putting the Cure into Strategic Procurement

## DENE STALK
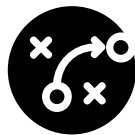
**Data Science Diploma**
BSc, MBA

Work Experience
Management Consulting, Internal Strategy, Tourism

## Importance of conducting a spend analysis

Identify cost saving opportunities

Facilitate strategic procurement decision making and sourcing
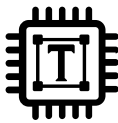
## Capstone Overview

Predict purchase order spend categories

6 weeks

Government of California's 2012-2015 purchase orders

NLP weighted voting ensemble using supplier name, item name and item description

Tools: Python, Pandas, Numpy, Sklearn, Tableau, MySQL

### Correctly Classified

**77% of line items**

**92% of Total Spend**

# Overview of Business Importance, Challenges and Machine Learning's value in procurement

## Why a Spend Analysis?

- Full spend visibility
- Identify cost saving opportunities
- Data-driven sourcing

## Challenges

- Disconnect across regions, languages and business units
- Labour intensive
- Prone to human error

## ML's Procurement Value

- Reduces need for manual intervention
- Timely strategic decision making
- Better inventory and supplier management

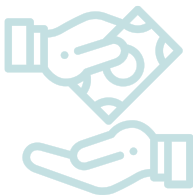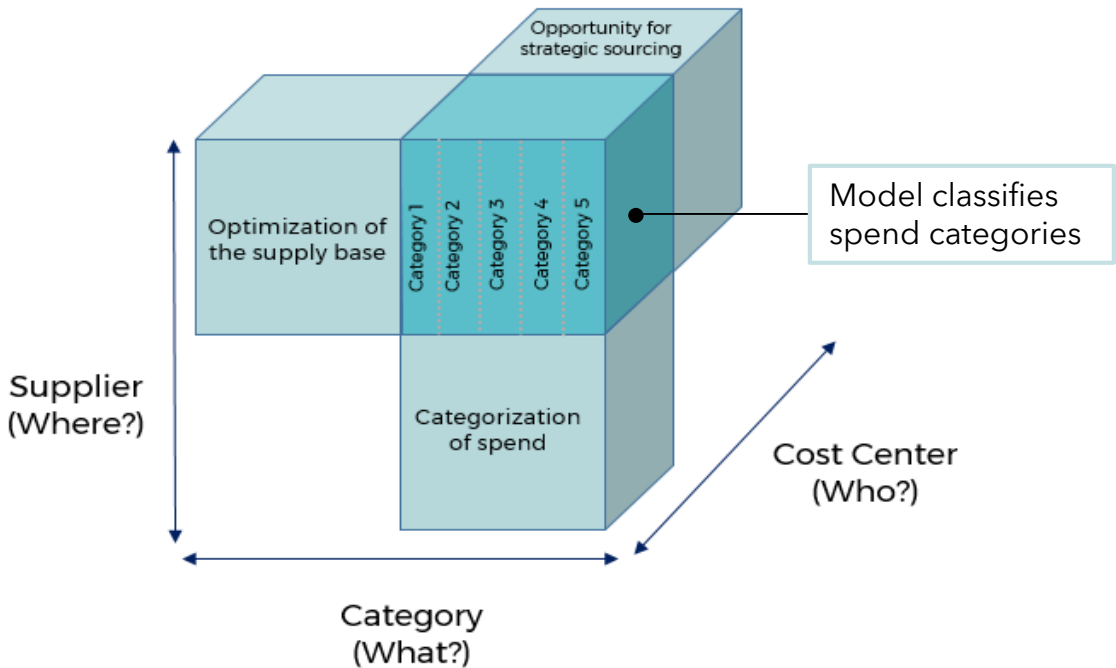# Importance of Correctly Classifying Spend Categories

Purchase Order (PO): document issued by a buyer committing to purchasing products or services

Takes an average of 1 year to identify and correct 100k misclassified line items

Increased spend visibility leads to lower Total Cost of Ownership[1]



Opportunity for strategic sourcing

Optimization of the supply base

Category 1
Category 2
Category 3
Category 4
Category 5

Model classifies spend categories

Supplier (Where?)

Categorization of spend

Cost Center (Who?)

Category (What?)

Improved buying power through strategic sourcing
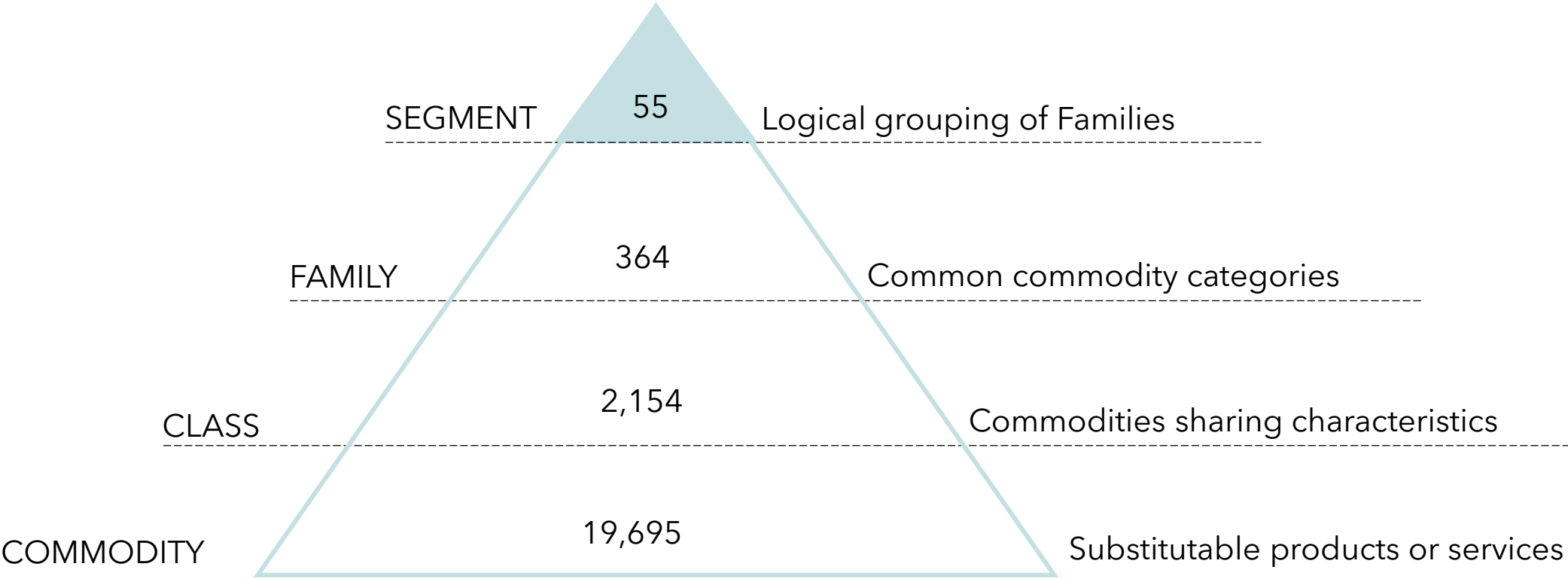
Optimizing inventory management processes

Managing risk exposure to maverick spending[2]

1 - Total Cost of Ownership is a combination of purchasing and carrying costs
2 - Maverick spending refers to purchases made outside of agreed contracts

# United Nation's Classification System

**UNSPSC** : United Nations Standard Products & Services Code

SEGMENT — 55 — Logical grouping of Families

FAMILY — 364 — Common commodity categories

CLASS — 2,154 — Commodities sharing characteristics

COMMODITY — 19,695 — Substitutable products or services

Source: Open Canada

# Data Collection and Description



**Purchase Order Data 2012-2015**
Government of California merged information from various procurement systems, then mapped the UNSPSC taxonomy

## Cleaned Dataset Description

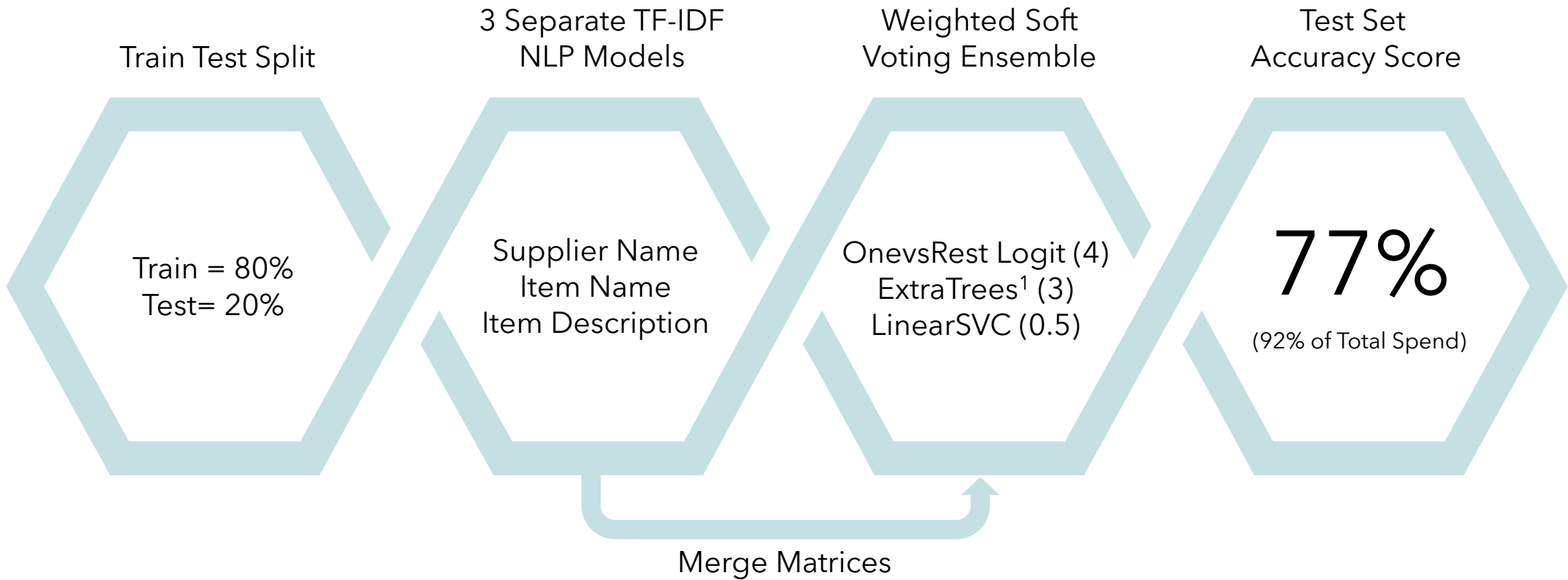| 3 fiscal years | 326k line items | 24k suppliers | 3.6m words |

| $147bn in Total Spend | 26 Spend categories[1] | | 171k unique item names 208k unique item descriptions |

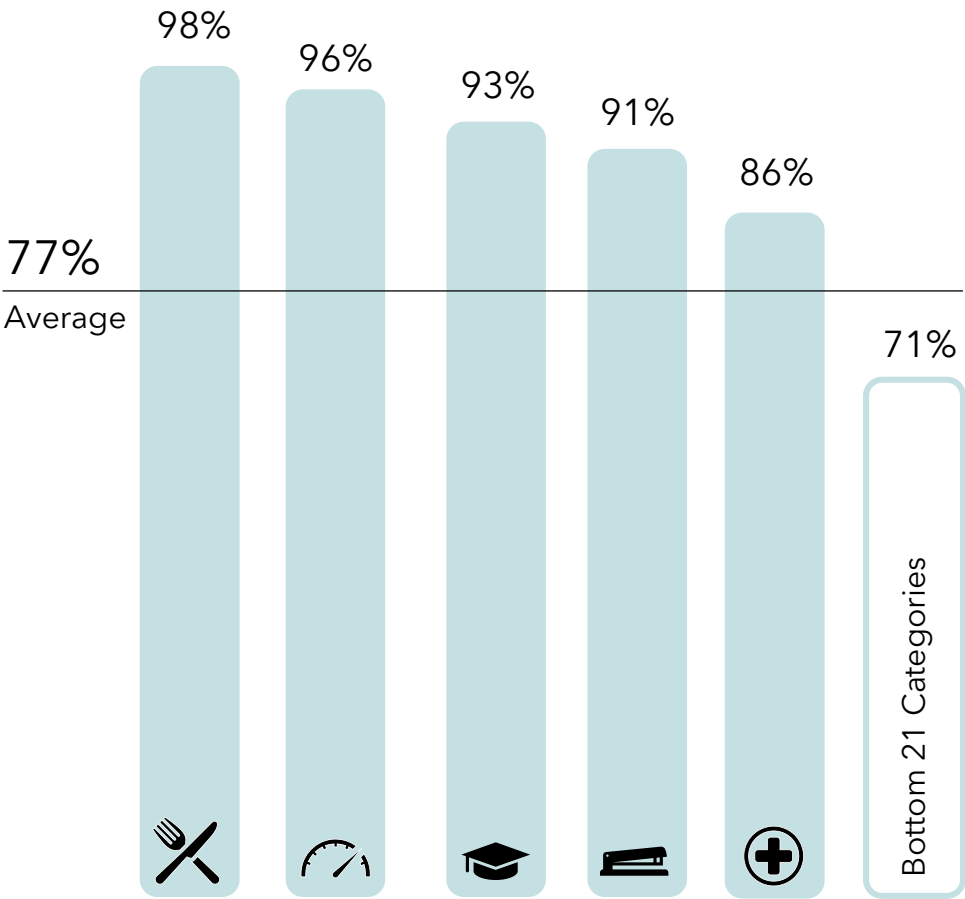1 – Bottom categories in terms of items lines were placed into either Other Services or Other Goods

Detailed data cleaning process available at github.com/denestalk

# Modelling Methodology

NLP weighted soft voting ensemble predicted 65k line items of unseen data with 77% accuracy

Train Test Split

3 Separate TF-IDF
NLP Models

Weighted Soft
Voting Ensemble

Test Set
Accuracy Score

Train = 80%
Test= 20%

Supplier Name
Item Name
Item Description

OnevsRest Logit (4)
ExtraTrees[1] (3)
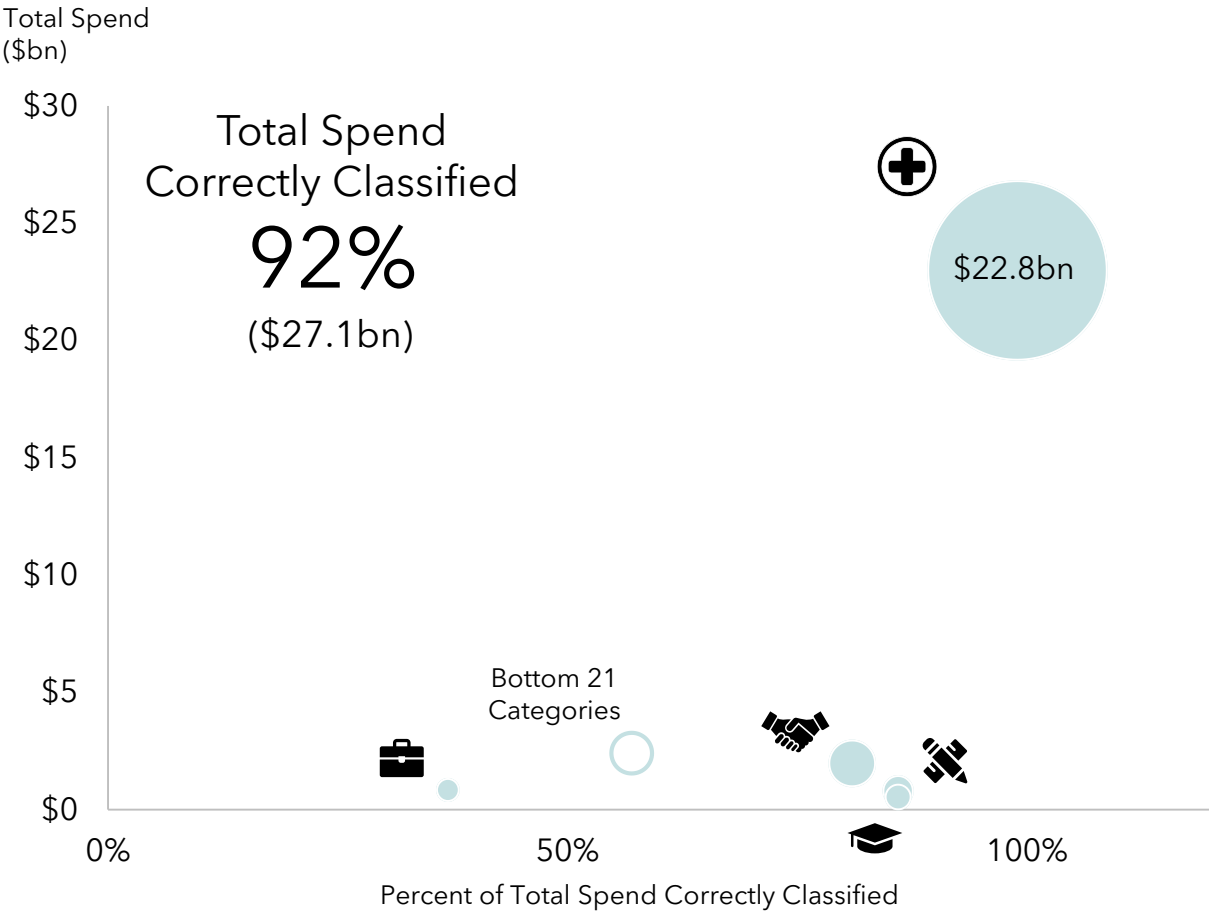LinearSVC (0.5)

77%

(92% of Total Spend)

Merge Matrices

1 - ExtraTrees differs from RandomForest by splitting nodes randomly and does not resample observations

# Results

## Top 5 Categories by Percent of Line Items Correctly Classified

98%
96%
93%
91%
86%

77%
Average

71%

Bottom 21 Categories

## Top 5 Categories by Total Spend
(bubble size indicates Total Spend correctly classified)

Total Spend
($bn)

$30

$25

$20

$15

$10

$5

$0

Total Spend
Correctly Classified
**92%**
($27.1bn)

$22.8bn

Bottom 21
Categories

0%          50%          100%

Percent of Total Spend Correctly Classified

⚒ Food & Beverage Products    ⊕ Health Services    🎓 Education & Training Services    🤝 Other Services

⌢ Fuel & Lubricants    ⌇ Office Supplies    💼 Professional Services    ⚒ Engineering Services

Category names have been abbreviated
Full breakdown available at github.com/denestalk

# Conclusion

## Model Strengths

- Relatively strong predictive value using only 3 free text columns

- Correctly classified the vast majority of spend

- 17 categories' line items correctly classified with over 70% accuracy

- Model is scalable

## Model Weaknesses

- Not tested on corporate data

- Health Services over represented in terms of Total Spend

- 3 categories < 50% of line items correctly predicted

## Further Analysis

- Meta-model incorporating quantity and unit price

- Deeper analysis on strategic value in relation to supplier management

- Test transferability on new Government of California data and corporate data

- Create an interactive dashboard

Category names have been abbreviated

# Thank you

**Dene Stalk**
denemstalk@gmail.com
416-573-5784

linkedin.com/in/denestalk
github.com/denestalk