# CE Diary Autocoder
## Demo

**Michell Li**

Data Science Intern

Civic Digital Fellow

August 6, 2019

# What is Machine Learning?

■ Definition: Machine learning is programming computers to **optimize a performance criterion** using example data or **past experience**

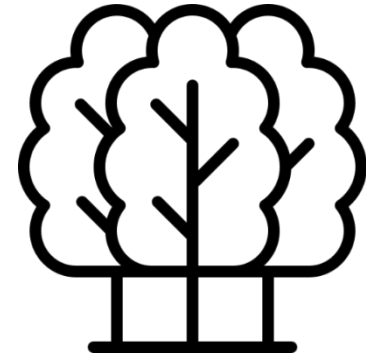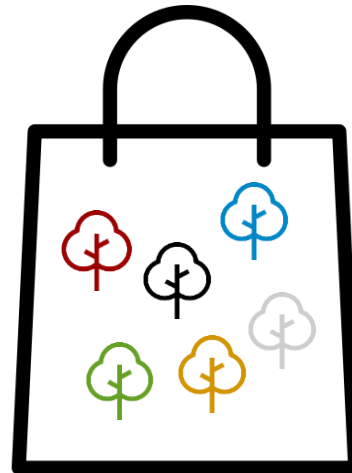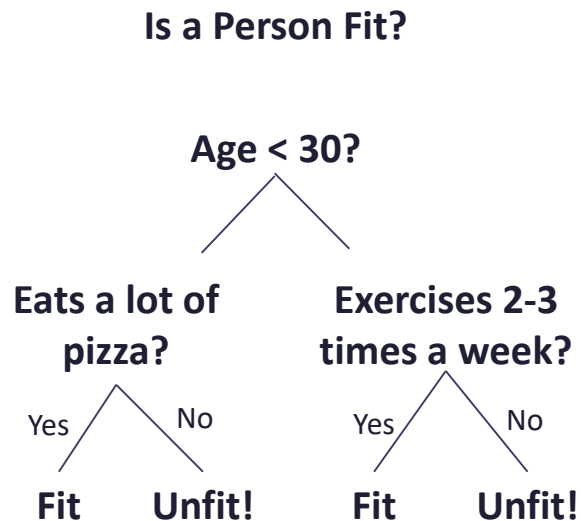■ Goal: Build a model that is a **good and useful approximation** to the data

## Machine Learning is Like Gardening

Gardener
You

Nutrients
Data

Seeds
Algorithm

Plants
Programs

# What is a Random Forest Model?

Random forest models are **bagged decision tree** models that split on a **subset of features** at each split.

### Decision Tree:

**Is a Person Fit?**

**Age < 30?**

**Eats a lot of pizza?**          **Exercises 2-3 times a week?**

Yes    No                          Yes    No

**Fit    Unfit!**                  **Fit    Unfit!**

# Motivation

The Bureau of Labor Statistics wants to automatically assign item codes in the Diary survey.

The process is currently labor intensive and expensive.

The existing autocoder is a rule based system that needs to account for special cases, which leads to inaccuracies.

Creating a new autocoder using machine learning can: **reduce costs, improve accuracy.**

# Contributions

**Models**

Four Models: ECLO, EFDB, EOTH, EMLS

**Data and Analysis**

Spell Checker created using Levenshtein, Jaro Winkler  and QWERTY distance

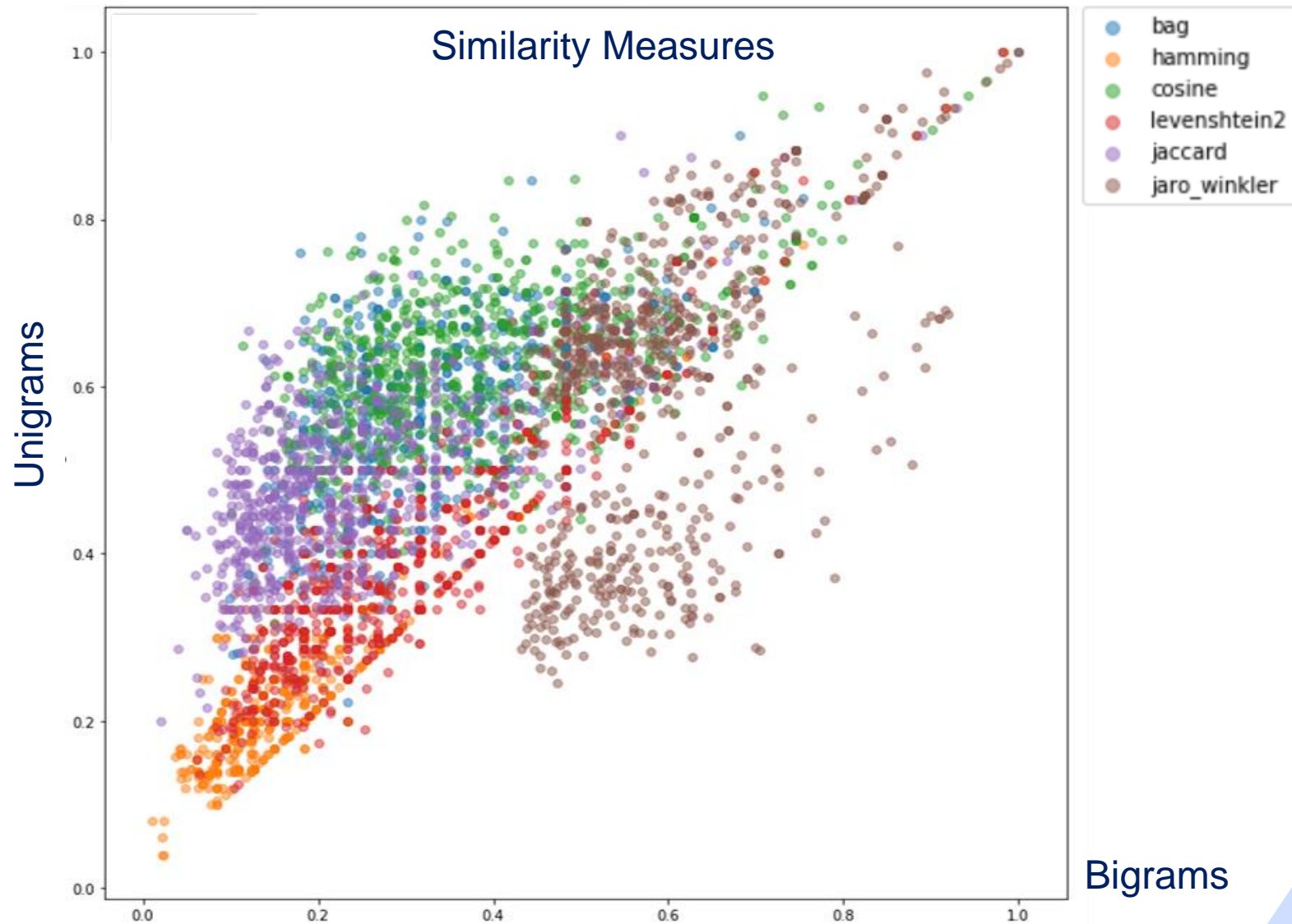Module for analyzing Unclassifiable diary items

Analytics Dashboard UI using Dash
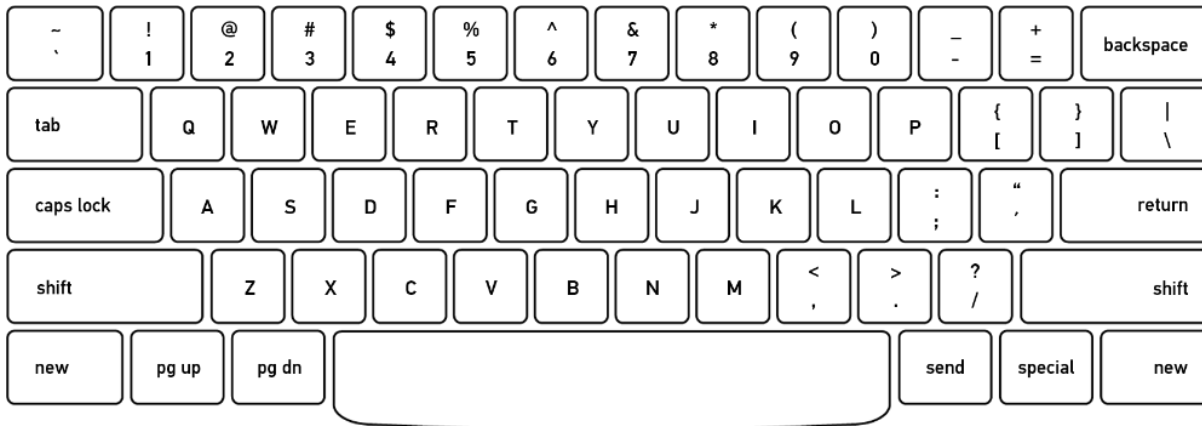
# Data Exploration

## # Unique Item Codes



**Total: 753**
(not including >900000)

# Dealing with Dirty Data



Similarity Measures

Unigrams

Bigrams

Legend:
- bag
- hamming
- cosine
- levenshtein2
- jaccard
- jaro_winkler

# Spell Checker

LABOT vs LAOBR to LABOR



Transposition/Replace
QWERTY Penalty:
T -> R = log(1)
B <-> O = log(4)

```
1   {
2       "spagh": "spaghetti",
3       "spagheti": "spaghetti",
4       "spaghettie": "spaghetti",
5       "spaghtti": "spaghetti",
6   }
```
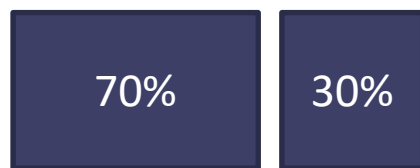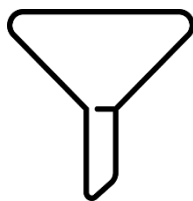
Insertion/Deletion
Penalty: 1

# Model Architecture

Item Descriptions are **spell checked and vectorized** before being split into training and testing data sets.
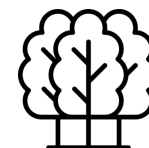
The data sets are then fed into respective Random Forest Models and **Item predictions are reflected on the dashboard.**

3 Years of SAS Data

| 70% | 30% |
|---|---|
| Train | Test |

Spell checked
TfidfVectorizer ngrams = (2,3)
shirt: sh, hi, ir, rt, shi, hir, irt

i.e.:
n_estimators = 200
max_depth = 100
criterion = entropy

# Train Model

Each RF model is trained using KFold cross validation.
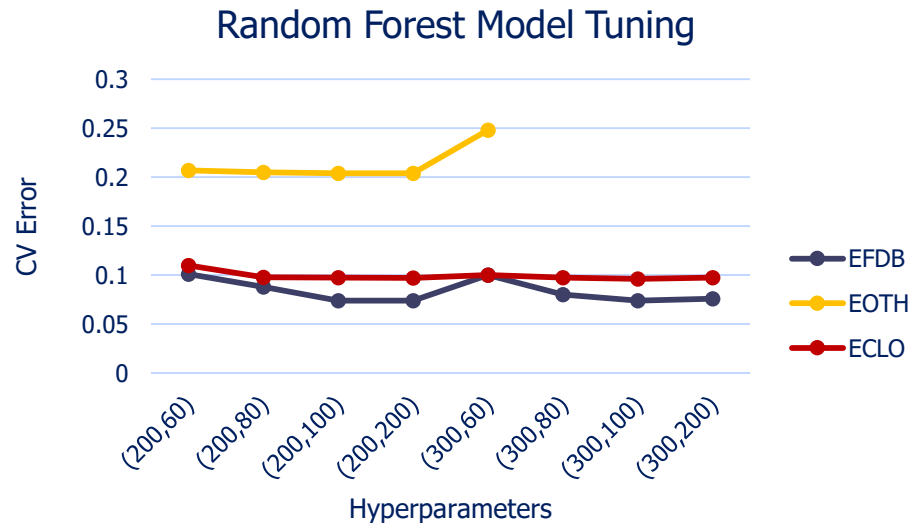
Hyperparameter settings:

n_estimators: [200, 300]
max_depth: [60, 80, 100, 200]
criterion: ['gini', 'entropy']

Total # of features

### Random Forest Model Tuning



Other models tried:
Logistic Regression, SVM, Decision Tree

# Demo

- Installation:
  - ▶ Anaconda/Python required
  - ▶ Pip install autocoder package
- Dashboard
- Download to Excel



```
Anaconda Prompt - python upload_component.py

(base) C:\>cd C:\Users\li_m\Documents\autocode\DiaryAutocoding\front

(base) C:\Users\li_m\Documents\autocode\DiaryAutocoding\front>python upload_component.py
Running on http://127.0.0.1:8050/
Debugger PIN: 770-477-768
 * Serving Flask app "upload_component" (lazy loading)
 * Environment: production
   WARNING: Do not use the development server in a production environment.
   Use a production WSGI server instead.
 * Debug mode: on
Running on http://127.0.0.1:8050/
Debugger PIN: 913-207-254
```
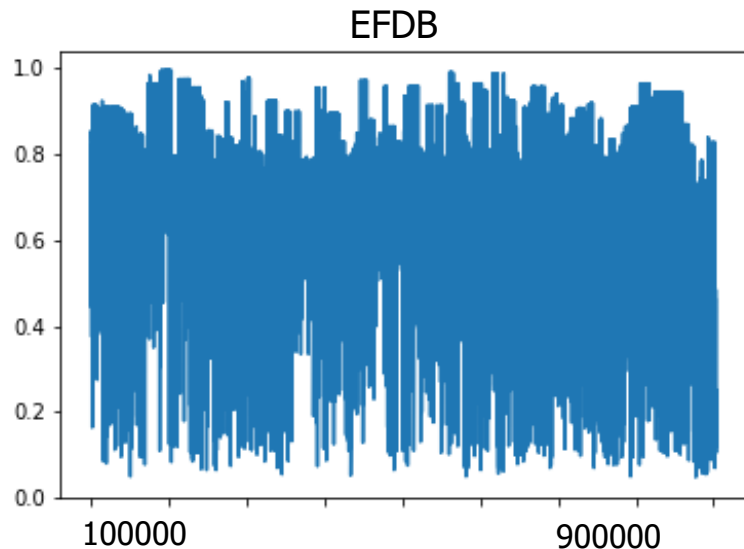
# Evaluation

| Model | Precision | Recall | F1 Score | Accuracy |
|---|---|---|---|---|
| Clothing | 0.91 | 0.9 | 0.9 | 0.9 |
| Food and Beverages | 0.95 | 0.94 | 0.94 | 0.94 |
| Meals | 0.98 | 0.98 | 0.98 | 0.98 |
| Other | 0.85 | 0.83 | 0.83 | 0.83 |

# Item Code Probability Distributions

EOTH

ECLO

EFDB

EOTH and EFDB show greater variability in predicted probabilities whereas ECLO shows greater confidence in its predictions.

# Evaluation



% Needed to Manually Code

Legend: ■ Current  ■ New

| Category | Current | New |
|---|---|---|
| EOTH | 45% | 16% |
| EMLS | 3% | 1% |
| EFDB | 28% | 6% |
| ECLO | 24% | 10% |

Distribution of Predicted Probabilities

Legend: ■ EFDB  ■ EOTH  ■ ECLO

# Threshold Evaluation

## Above Threshold

| Model | Precision | Recall | F1 Score | Accuracy |
|-------|-----------|--------|----------|----------|
| Clothing | 0.97 | 0.96 | 0.97 | 0.97 |
| Food and Beverages | 0.99 | 0.99 | 0.99 | 0.99 |
| Other | 0.97 | 0.97 | 0.97 | 0.97 |

## Below Threshold

| Model | Precision | Recall | F1 Score | Accuracy |
|-------|-----------|--------|----------|----------|
| Clothing | 0.802 | 0.7 | 0.71 | 0.7 |
| Food and Beverages | 0.87 | 0.85 | 0.85 | 0.85 |
| Other | 0.79 | 0.73 | 0.72 | 0.73 |

# Takeaways

1. Even though **EFDB** has the most number of unique item codes (185), **the codes are the most separable** (least overlap); **ECLO** item codes have **the most overlap**

2. Thresholds can be lowered (less manual coding) to a certain extent without accuracy compromise

3. Misspellings does not impact accuracy that much
   - (ECLO accuracy +0.4% after misspellings fixed)

# Future Work

This autocoder is a **promising proof of concept** for use of machine learning at the BLS.

- **Moving away from Census NPC coding:** Semi-supervised learning methods can predict item codes using past data's target variable without the need for its own

- **Reducing model error:** boosting methods, fix misspellings, dimensionality reduction, more training data (more CPUs), include store name

- **Thresholding for manual classification:** One vs. Rest Classifier ROC curve to determine specific threshold

- **Incorporating BLS expertise:** Diary specific vector embeddings using Genism's Word2Vec

# Thank You!

Questions?

# Appendix

# Ethics of Machine Learning

Machine Learning uses **past data sets to predict outcomes**, meaning a model is only as good as its data.

**Machine Learning in the Government**

Bias and discrimination: Government Crime Classification Tool Racially Biased
Erosion of Privacy: Chinese Government Launches Social Credit System, NYC Patrolling Officers Wear Body Cameras

**Machine Learning at the BLS**

Models have **preferences**.

Automation contributes **to workforce displacement**.

Automation automates **human biases**.

```
"pepperoi": "pepper",
```

| ItemType | ItemCode | ItemDescription |
|----------|----------|-----------------|
| EOTH | 821132 | ACRYLIC NAILS |
| EOTH | 870048 | NAILS |
| EOTH | 314120 | COMMON NAILS |

*Under what circumstances is autocoding worthwhile?*