# Total Recall

Samuel Boardman, Khola Jamshad, Riku Kurama, Yucong Lei, Shivani Prabala

# The Problem

**Motivation**

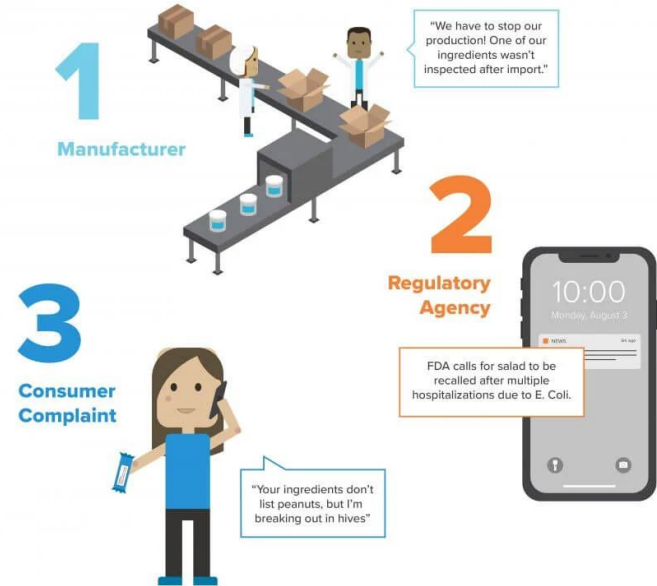FDA to suspend quality-control program for food testing due to staff cuts

**Research Question**

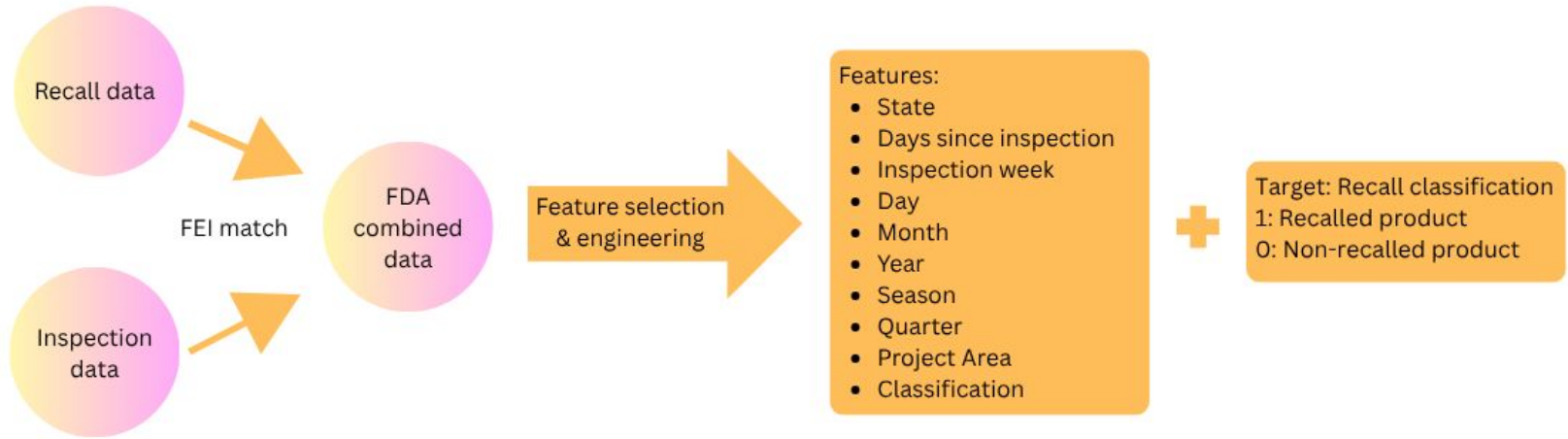Can we predict food products likely to be recalled by the FDA?

**Intended Impact**
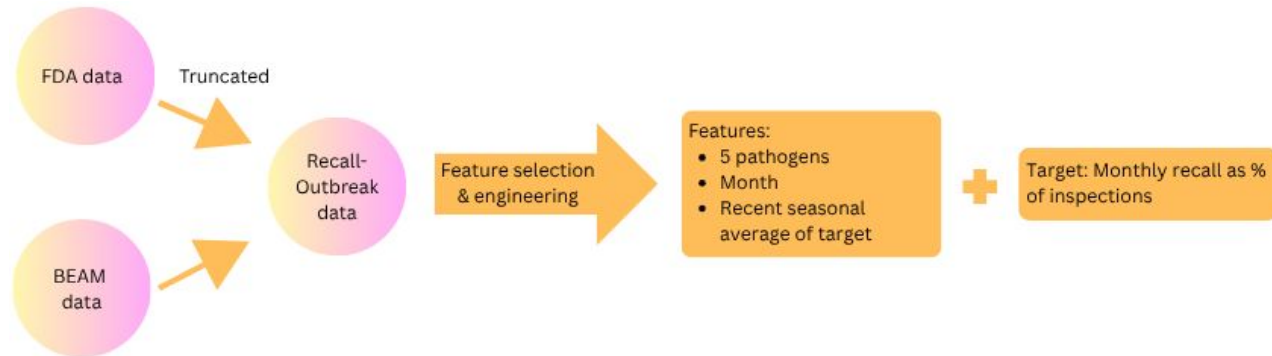
Reduced contaminated products reaching consumers

Reduced costs to manufacturers

# Datasets



Recall data → (FEI match) → FDA combined data → Feature selection & engineering →

Inspection data →

Features:
- State
- Days since inspection
- Inspection week
- Day
- Month
- Year
- Season
- Quarter
- Project Area
- Classification

Target: Recall classification
1: Recalled product
0: Non-recalled product

**Future Work setup:**

FDA data → (Truncated) → Recall-Outbreak data → Feature selection & engineering →

BEAM data →

Features:
- 5 pathogens
- Month
- Recent seasonal average of target

Target: Monthly recall as % of inspections
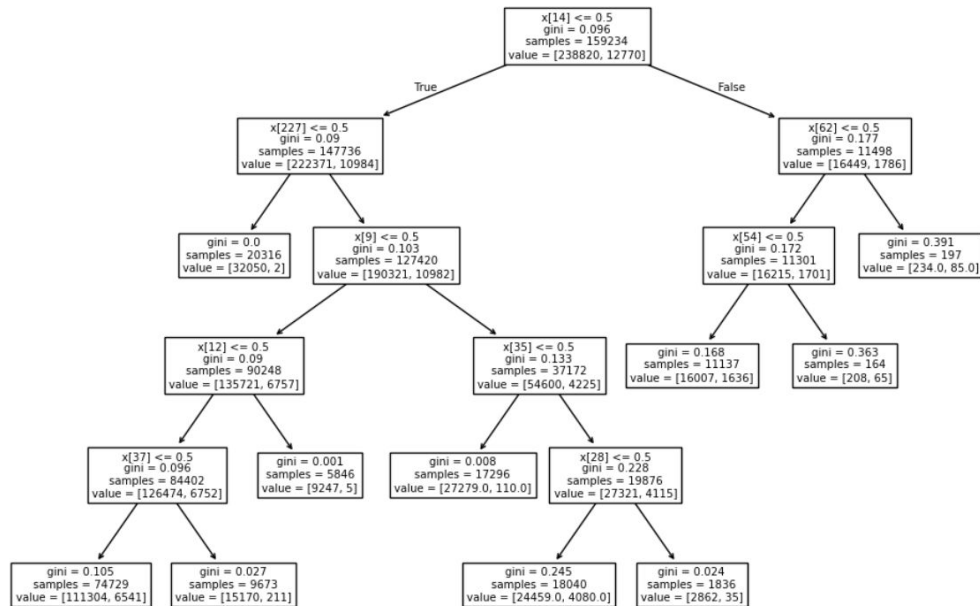
# Random Forest Classification

- Model overview
  a. Use of thresholds
- Hyperparameter tuning
  a. Evaluation metric
- Final hyperparameter choice
  a. n_estimators = 100
  b. max_depth = 10
  c. class_weight = 'balanced'

# Support Vector Machine Classification

| SVM with class_weight = 'balanced' | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| | 0 | 0.92 | 0.57 | 0.70 | 29779 |
| | 1 | 0.09 | 0.48 | 0.15 | 2676 |
| | accuracy | | | 0.56 | 32455 |
| | macro avg | 0.51 | 0.52 | 0.43 | 32455 |
| | weighted avg | 0.86 | 0.56 | 0.66 | 32455 |

**SVM with class_weight = 'balanced'; grid search for optimal parameters**

```
Fitting 5 folds for each of 12 candidates, totalling 60 fits
✅ Best Parameters: {'C': 0.1, 'class_weight': 'balanced', 'gamma': 'scale', 'kernel': 'rbf'}
✅ Best Recall Score (CV average): 0.5474658901592778
```

Test Classification Report:

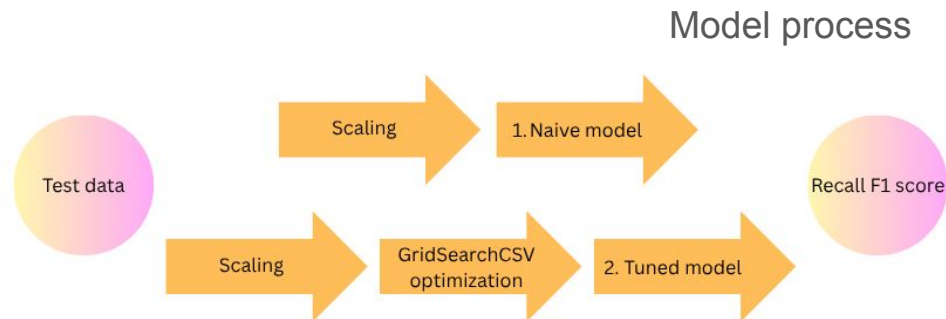| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.92 | 0.52 | 0.66 | 29779 |
| 1 | 0.09 | 0.53 | 0.15 | 2676 |
| accuracy | | | 0.52 | 32455 |
| macro avg | 0.51 | 0.52 | 0.41 | 32455 |
| weighted avg | 0.86 | 0.52 | 0.62 | 32455 |

**SVM with class_weight = 'balanced'; grid search for optimal parameters on trimmed dataset**

Classification Report:

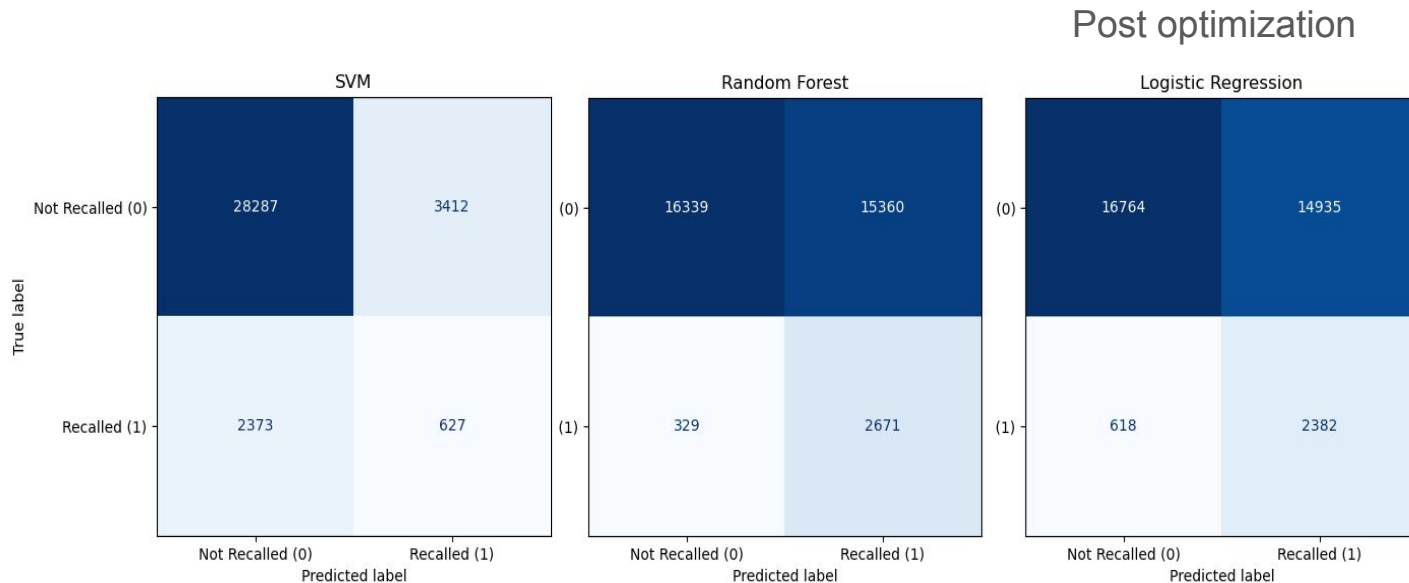| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.22 | 0.58 | 0.32 | 594 |
| 1 | 0.88 | 0.60 | 0.71 | 3064 |
| accuracy | | | 0.60 | 3658 |
| macro avg | 0.55 | 0.59 | 0.51 | 3658 |
| weighted avg | 0.77 | 0.60 | 0.65 | 3658 |

# Model Selection

- Dummy (majority class) as baseline

- LR and RF F1 score improved by ~ 25%

- SVM still struggled to predict recalls



Test data → Scaling → 1. Naive model → Recall F1 score

Scaling → GridSearchCSV optimization → 2. Tuned model

## Conclusions:

Post optimization

- **Random forest** is our **chosen model** for its ability to predict recalls.

- **Predicting recalls is difficult** with current data especially due to imbalance.



**SVM**

|  | Not Recalled (0) | Recalled (1) |
|---|---|---|
| Not Recalled (0) | 28287 | 3412 |
| Recalled (1) | 2373 | 627 |

**Random Forest**

|  | Not Recalled (0) | Recalled (1) |
|---|---|---|
| (0) | 16339 | 15360 |
| (1) | 329 | 2671 |

**Logistic Regression**

|  | Not Recalled (0) | Recalled (1) |
|---|---|---|
| (0) | 16764 | 14935 |
| (1) | 618 | 2382 |

# Future Work

Predicting recall percentages in future months using a VAR (vector autoregression) model.

**Variables used:** recall percentage, numbers of outbreaks with 5 types of pathogens, recent seasonal average of recall percentage

Our current model only performs roughly as well as a baseline seasonal model.

**Potential improvement:**
Including more spatial information, climate reasons for recall, and by adding regularization.



Recall percentage by month