
Predicting Chicago Food Inspection Failures

Jingzhi Yang Lingyue Ji Nate Assefa

1. Data Description

1.1. Objective

Our research question is clear and focused: *Can we predict, at the moment an inspection occurs, whether a Chicago food facility will fail?* This question specifies both the prediction target, which is inspection failure, and the timing constraint, that only pre-inspection information may be used.

The City of Chicago Food Inspections dataset is well suited to this goal. It provides detailed facility characteristics, inspection metadata, location, and free-text violation notes. These fields capture operational risk factors and historical outcomes that are essential for modeling future failures.

Our research strategy is a supervised machine learning approach. We first clean and normalize the raw data, engineer features such as an ordinal risk score and calendar variables, and then train classification models using the cleaned training set while holding out a test set for final evaluation. This strategy matches the research question because it builds a predictive model while respecting the constraint that only information available before the inspection is included. The public documentation and data dictionary of the dataset serve as the primary codebook and reference for variable definitions.

1.2. Team Collaboration

Lingyue found the raw data, exported cleaned training and testing CSV files, and created the repository. Nate cleaned the data, checked schema, and performed visualization. Jingzhi managed the files, created the data dictionary, and wrote the report.

1.3. Variables

The dataset contains both original inspection fields and engineered features for modeling. Key original variables include:

- **inspection_id**: Unique identifier for each inspection event
- **inspection_date**: Calendar date of the inspection; used for time-based analysis
- **facility_type**: Establishment category (e.g., restaurant,

grocery store, school) capturing operational risk

- **risk**: City-defined risk level with three categories: Risk 1 (High), Risk 2 (Medium), Risk 3 (Low)
- **inspection_type**: Reason for inspection (e.g., canvass, complaint, license, re-inspection)
- **zip, city, latitude, longitude**: Facility location variables for geographic analysis
- **results**: Official outcome string (pass, pass w/ conditions, fail, or out of business); defines the target label.
- **violations**: Free-text list of cited health-code violations, often semi-structured with numbered items

We also create engineered features to support modeling:

- **viol_count**: Count of numbered violations parsed from the free-text field
- **risk_ord**: Ordinal encoding of risk (High=3, Medium=2, Low=1)
- **ins_year, ins_month, ins_wday**: Year, month, and weekday extracted from the inspection date
- **target_fail**: Binary target variable; equals 1 if the inspection result is “fail,” 0 otherwise

These variables together describe the facility, its risk context, the inspection circumstances, and outcomes. They form the basis for predicting inspection failure while avoiding data leakage by excluding violation text from the predictive model.

1.4. Data Cleaning

Each data split contains 1,000 inspections. After cleaning, the tables hold 18 columns, and about 19% of training records show a failed inspection. Facility type has about 41 categories, with restaurants most common. Inspection type has about 19 categories, and risk has three. Coordinates cluster tightly around Chicago (about 41.88° N, −87.67° W). Years range from 2010 to 2025 with even coverage. Failure rates rise with higher risk, though risk 2 shows a slightly higher rate than risk 3.

Cleaning follows a single reproducible pipeline. We drop duplicate inspection IDs and standardize text. City names convert to title case, and a new field flags Chicago locations. ZIP codes convert from floats to five-digit strings with leading zeros. Dates parse to `datetime` objects, and we derive year, month, and weekday. Result strings normalize to a compact set (“pass,” “pass w/ conditions,” “fail,” and similar). We encode the risk field to an ordinal integer. Violation notes become a numeric count using a layered parsing strategy. Latitude and longitude outside conservative Chicago bounds become missing.

1.5. Challenges

The first dataset we considered was from Kaggle and was eight years ago. After discussion, we switched to the current dataset, which is more up-to-date.

Text fields contained inconsistent casing and spelling, so we standardized them. Some cells contained the string “nan,” which we converted to missing. ZIP codes required conversion from float to string. Violation text is semi-structured, so we used a conservative counting method. The test set includes facility types not present in training, so we map unseen categories to an “Other” bucket and use encoders that tolerate unknown levels.

2. Pre-Analysis Plan

2.1. Method Overview

We treat this project as a supervised binary classification task. The goal is to predict whether a food inspection will *fail* (`target_fail = 1`) using only information known before the inspection begins.

Our team builds a clear and repeatable pipeline. The pipeline first validates the dataset schema and cleans both the training and testing data in the same way. It then engineers new features while keeping the time constraint intact. After that, it benchmarks several models, starting from simple interpretable baselines and moving to more complex tree-based ensembles.

We train and tune all models using stratified cross-validation within the training set. The test set remains untouched until the final evaluation to measure how well the model generalizes. All steps run in a single reproducible script. Each step is logged to ensure transparency and auditability.

2.2. Team Collaboration

Lingyue prepared the cleaned training and testing files and verified schema consistency. Nate implemented the modeling pipeline, trained the baseline and ensemble models, and maintained the experiment structure. Jingzhi evaluated

the validation results, checked the reproducibility of the pipeline, and drafted the written report.

2.3. Data Cleaning

We apply a consistent cleaning process to both the training and testing datasets. This prevents bias between the two splits. Our code first checks that all required columns exist and removes any duplicate `inspection_id`.

We standardize all text fields. Each string is trimmed, converted to lowercase or title case, and mapped into consistent categories. The `results` field is reduced to a small, clear set of outcomes such as “pass,” “fail,” and “out of business.”

City names are normalized, and a new flag marks whether a facility is located in Chicago. ZIP codes are converted to five-digit strings, and invalid or missing ones are set to NA.

We parse the `inspection_date` column into a proper date format. From it, we extract the inspection year, month, and weekday.

Risk labels are encoded as ordered integers, where high risk receives the largest value. Violation notes are parsed from semi-structured text to count the number of cited violations. When the field contains non-informative placeholder strings, the parser applies a conservative rule and assigns a count of 1. When the field is truly missing (NA), the count is 0. We keep only this numeric count to avoid text leakage.

Finally, we check latitude and longitude values. Points outside the reasonable Chicago range (41.60–42.10° N, 87.95–87.50° W) are marked as missing.

After cleaning, each dataset contains 18 columns. Restaurants remain the most common facility type. The data span 2010–2025 and cover all three risk levels and 19 inspection types.

2.4. Modeling Plan and Justification

We use a small but diverse set of models that work well with mixed tabular data. Each model serves a clear purpose in our benchmarking process.

1. **Logistic Regression:** This model offers a simple and interpretable starting point. It produces calibrated probabilities and helps us measure how much more complex models improve performance. We encode categorical features and apply ℓ_2 regularization to prevent overfitting. The model estimates the probability of failure as:

$$\Pr(y=1 | x) = \sigma \left(\beta_0 + \sum_j \beta_j x_j \right),$$

where σ is the logistic function.

2. **Decision Tree:** This model captures simple feature interactions and threshold effects. It produces decision rules that are easy to interpret.
3. **Random Forest:** This model combines many decision trees to reduce variance. It is more stable than a single tree and handles outliers and correlated features effectively.
4. **Gradient-Boosted Trees:** These models often achieve the best performance on structured data. They build trees sequentially, learning from previous errors, and can model complex nonlinear relationships with built-in regularization.

We divide our features into three main types. Continuous variables include latitude, longitude, `viol_count`, and calendar fields such as year and month. The ordinal variable is `risk_ord`, which represents the inspection risk level. Categorical variables include `facility_type`, `inspection_type`, ZIP code, and `city_is_chicago`. We use one-hot encoding and handle unseen categories safely.

We do not use the raw `violations` text to avoid information leakage. Future versions of the model may use text features with careful cross-validation.

2.5. Training, Tuning, and Evaluation

We keep the provided test set completely separate from model development. The test data are used only once, at the very end, to measure real-world performance. All training, tuning, and model selection happen inside the training set.

We train each model using stratified k -fold cross-validation with $k=5$. This method keeps the same fail/pass ratio in every fold. It helps us make better use of limited data and produces stable performance estimates.

We tune model hyperparameters with either randomized search or grid search. The search covers key parameters such as tree depth, learning rate, minimum child weight, and subsampling rate. These ranges are chosen based on prior experience with tabular datasets.

Inspection failures are less common but more important to predict. Because of this imbalance, we report both threshold-free and threshold-based metrics. For threshold-free evaluation, we use ROC-AUC and PR-AUC to measure ranking quality. For threshold-based evaluation, we report Recall, Precision, and F_1 for the fail class. We also check model calibration with reliability curves. These plots show how well predicted probabilities match observed outcomes. To support interpretability, we will visualize feature importances and SHAP values.

We address class imbalance in two main ways. First,

we apply class-weighted loss functions, such as `class_weight='balanced'` in logistic models or `scale_pos_weight` in gradient boosting. Second, we adjust classification thresholds on validation folds to balance recall and precision. We avoid oversampling within folds unless it becomes strictly necessary. This prevents data leakage and keeps evaluation honest.

2.6. Validation Plan

We select models and thresholds by stratified 5-fold CV on the training set. We then run a single, final evaluation on the held-out test set. A gap between CV and test indicates overfitting or drift; in that case we will tighten regularization, simplify features, or adopt temporal CV.

2.7. Implementation and Reproducibility

We implement the full pipeline as a single script with fixed seeds and logged configs. We version the cleaning function and apply it identically to train/test. We export cleaned CSVs, a meta-schema JSON, and EDA plots. Unknown categorical levels are routed to an explicit “unknown” bucket.

3. Results

3.1. Team Collaboration

Before running the models, we met to align the implementation with the pre-analysis plan and divide responsibilities. Lingyue first verified that the cleaned datasets and the engineered features met the schema expected by the modeling code, ensuring that the train-test transformations were applied consistently. Nate built the modeling pipeline, implemented the cross-validation and training procedures, and produced the resulting model outputs and plots. And Jingzhi evaluated these results, compared cross-validation performance with test-set metrics, and synthesized the findings in the written report.

3.2. Descriptive Patterns

We begin by examining the cleaned training sample of 1,000 inspections. The failure rate is about 19%, so the positive class (fail) is much smaller than the negative class. Histograms of key predictors (Figures 1, 2, 3, 4, 5, 6, and 7) show clear patterns. The ordinal risk score is skewed toward high risk, with most facilities labeled as Risk 1. The violation count is right-skewed as well, with many inspections reporting only one or two violations and a long tail of cases with ten or more violations.

The spatial coordinates cluster tightly around the Chicago city center. The inspection year, month, and weekday variables appear roughly balanced over the 2010–2025 period, which suggests good temporal coverage. Failure rates also

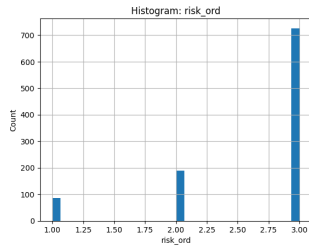


Figure 1. Histogram of the ordinal risk score.

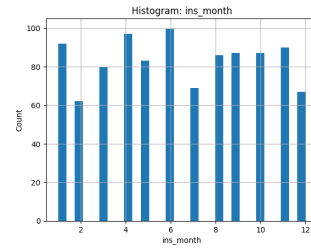


Figure 6. Distribution of inspection months.

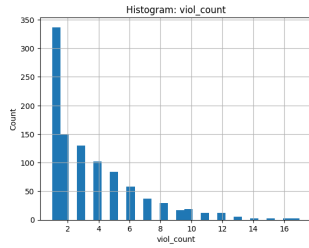


Figure 2. Histogram of violation counts.

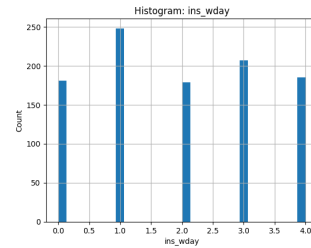


Figure 7. Distribution of inspection weekdays.

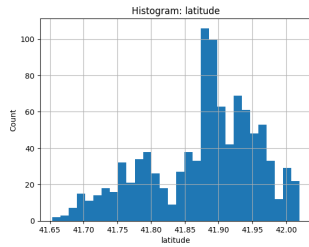


Figure 3. Distribution of facility latitude.

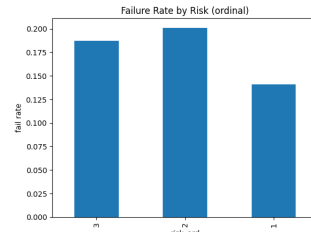


Figure 8. Failure rate by risk level.

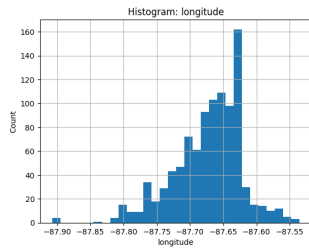


Figure 4. Distribution of facility longitude.

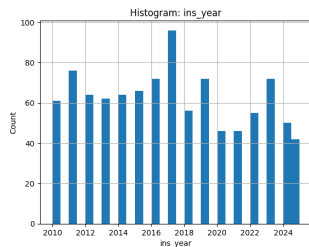


Figure 5. Distribution of inspection years.

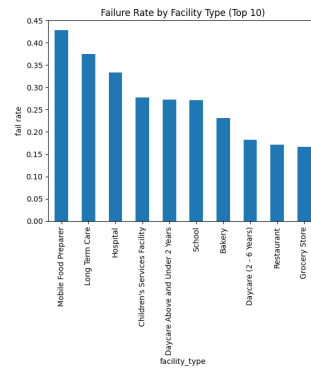


Figure 9. Failure rate by facility type (top 10).

vary across risk levels and facility types. Higher-risk establishments fail more often, and medium-risk sites fail slightly more often than low-risk sites (Figure 8). Among the ten most common facility types, mobile food dispensers and long-term care facilities show the highest failure rates. Grocery stores and several child-care related facilities fall in the middle of the distribution (Figure 9). Together, these plots show that the pre-inspection metadata contain useful information for predicting inspection outcomes.

3.3. Model Implementation and Training

Using the cleaned training data, we treat failure prediction as a supervised binary classification task. The feature matrix includes all pre-inspection variables in our schema, such as facility type, risk category, inspection program type, ZIP code, city flag, latitude and longitude, calendar features, and the numeric violation count. We follow the no-leakage rule from the pre-analysis plan and exclude free-text violation descriptions and the raw results string.

We implement a single preprocessing-and-model pipeline using `scikit-learn`. The pipeline encodes categorical variables with a `OneHotEncoder` that safely handles unseen levels. It passes numeric variables through unchanged, which creates a sparse design matrix suited to both linear and tree-based models. On top of this shared transformer, we fit four classifiers: logistic regression with ℓ_2 regularization and class-weighted loss, a single decision tree with class-weighted impurity, a random forest with 200 trees and balanced class weights, and a gradient-boosted decision tree model.

We select models using stratified 5-fold cross-validation on the training set. For each model, we compute ROC-AUC, PR-AUC (average precision), precision, recall, and F_1 for the fail class. We keep most hyperparameters at reasonable defaults and apply class-weighting in the linear and tree models to help address class imbalance.

3.4. Cross-Validation Benchmark

Table 1 reports the mean cross-validation performance for all four models. Logistic regression shows the strongest overall ranking quality, with ROC-AUC 0.82 and PR-AUC 0.58. It also shows a balanced trade-off between precision and recall. Gradient boosting performs well in ROC-AUC (0.80) but places more weight on precision than recall, which causes the model to miss more failures. The random forest model produces very high precision but extremely low recall and correctly identifies only about 8% of failed inspections. The single decision tree acts as a simple baseline and has the lowest ROC-AUC value (0.63).

Given its strong ROC-AUC, solid PR-AUC, and relatively high F_1 , we select the logistic regression model as the pri-

Table 1. Mean 5-fold cross-validation performance on the training set.

Model	ROC-AUC	PR-AUC	Precision	Recall	F_1
Logistic regression	0.824	0.582	0.402	0.651	0.497
Gradient boosting	0.802	0.514	0.657	0.269	0.380
Random forest	0.792	0.487	0.724	0.076	0.134
Decision tree	0.626	0.273	0.382	0.404	0.391

mary model for test-set evaluation. The results show that a linear classifier with one-hot encoded features can outperform more flexible ensemble methods in this setting. We initially expected gradient boosting to lead the benchmark, given its usual advantage on tabular data, so the strength of logistic regression served as a useful sanity check on our feature engineering. Together, these patterns suggest that most of the predictive structure in our features can be captured by approximately linear boundaries at this sample size.

3.5. Test-Set Performance

We refit the chosen logistic regression pipeline on the full training set and evaluate it once on the held-out test set of 1,000 inspections. The model shows strong discrimination on this new data:

- ROC-AUC = 0.811,
- PR-AUC = 0.548,
- Precision = 0.429,
- Recall = 0.691,
- F_1 = 0.530.

The ROC curve in Figure 10 stays well above the diagonal baseline. The PR curve in Figure 11 also lies well above the horizontal line for the base failure rate of about 19%. These patterns show that the model can rank inspections in a useful way. Higher predicted scores correspond to higher observed failure rates, which means the model can help prioritize high-risk inspections.

At the default probability threshold of 0.5, the confusion matrix in Figure 12 reports 628 true negatives, 178 false positives, 134 true positives, and 60 false negatives. These numbers mean that the model identifies about 69% of inspections that actually fail. They also show that about 43% of predicted failures are correct. This operating point favors recall for the rare fail class, but it also produces extra follow-up for some inspections that would have passed.

3.6. Initial Result Analysis

These results show that a simple model can still be useful. The logistic regression model uses only pre-inspection meta-

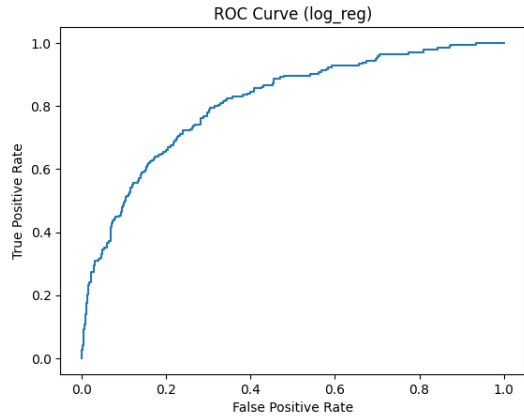


Figure 10. ROC curve for the logistic regression model.

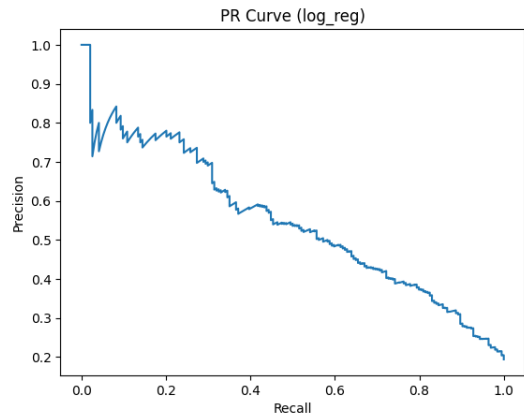


Figure 11. Precision–Recall curve for the logistic regression model.

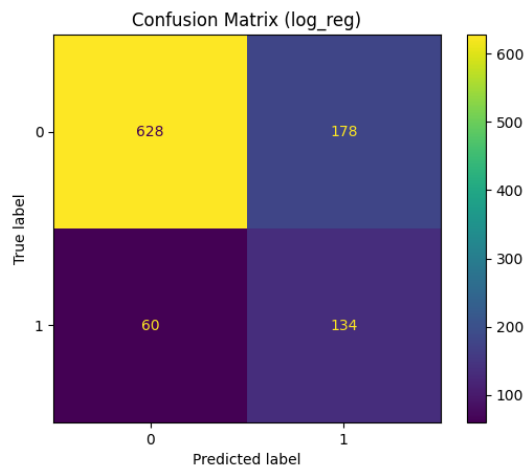


Figure 12. Confusion matrix for the logistic regression model on the test set.

data and a few engineered features, yet it predicts failures with reasonable accuracy. The small gap between cross-validation and test performance (e.g., ROC–AUC 0.82 vs. 0.81) suggests that the model does not overfit badly. The model also appears to generalize well to new inspections drawn from the same process.

From a policy perspective, a model with recall around 70% could help cities prioritize limited inspection resources. Inspectors could use predicted failure probabilities to decide which facilities to visit first. They could also assign more experienced staff to inspections with higher predicted risk. However, a precision of about 43% means that many flagged inspections would still pass. This trade-off may or may not be acceptable, depending on the cost of extra inspections and the harm caused by missed failures.

The comparison across models shows that more complex tree-based ensembles do not always perform better in this setting. Instead, careful handling of class imbalance, thoughtful feature engineering, and a well-calibrated linear model already capture most of the signal in the data. Because our sample size is modest and inspections span 2010–2025, future work should check whether performance degrades for the most recent years, especially if inspection policies have changed. They could adjust decision thresholds for different operational goals or add temporal validation to account for policy or behavior changes over time. Moreover, it could explore text features from violation histories in a way that still respects the no-leakage constraint.

References

City of Chicago. Food Inspections: City of Chicago Data Portal. <https://data.cityofchicago.org/Health-Human-Services/Food-Inspections/4ijn-s7e5>, 2025. Accessed September 26, 2025.