

BUSINESS SURVIVAL ANALYSIS SAN FRANCISCO DISTRICT 3

Presented by Abirami

Presentation Held on May 8th 2025

BACKGROUND

Context:

- District 3 includes diverse neighborhoods like Chinatown, North Beach, and parts of the Financial District, each with a unique business landscape.

Executive Summary:

- Goal: Predict administrative closure of businesses in District 3
- Tool: Survival Analysis with Cox Proportional Hazards model
- Outcome: Insights to inform targeted funding decisions

TARGET AUDIENCE

San Francisco District 3 funding team

GOAL

- The Main Goal of this project is to analyze the factors influencing the survival of businesses in San Francisco and for Supervisor District 3. Specifically, aim to predict the likelihood of a business being administratively closed over time using survival analysis techniques.

PROBLEM STATEMENT

Why Model Business Survival in District 3?

- Targeted Resource Allocation
- Informed Funding Strategies
- Performance Measurement
- Prioritization of Investment

DATA COLLECTION

Data Source

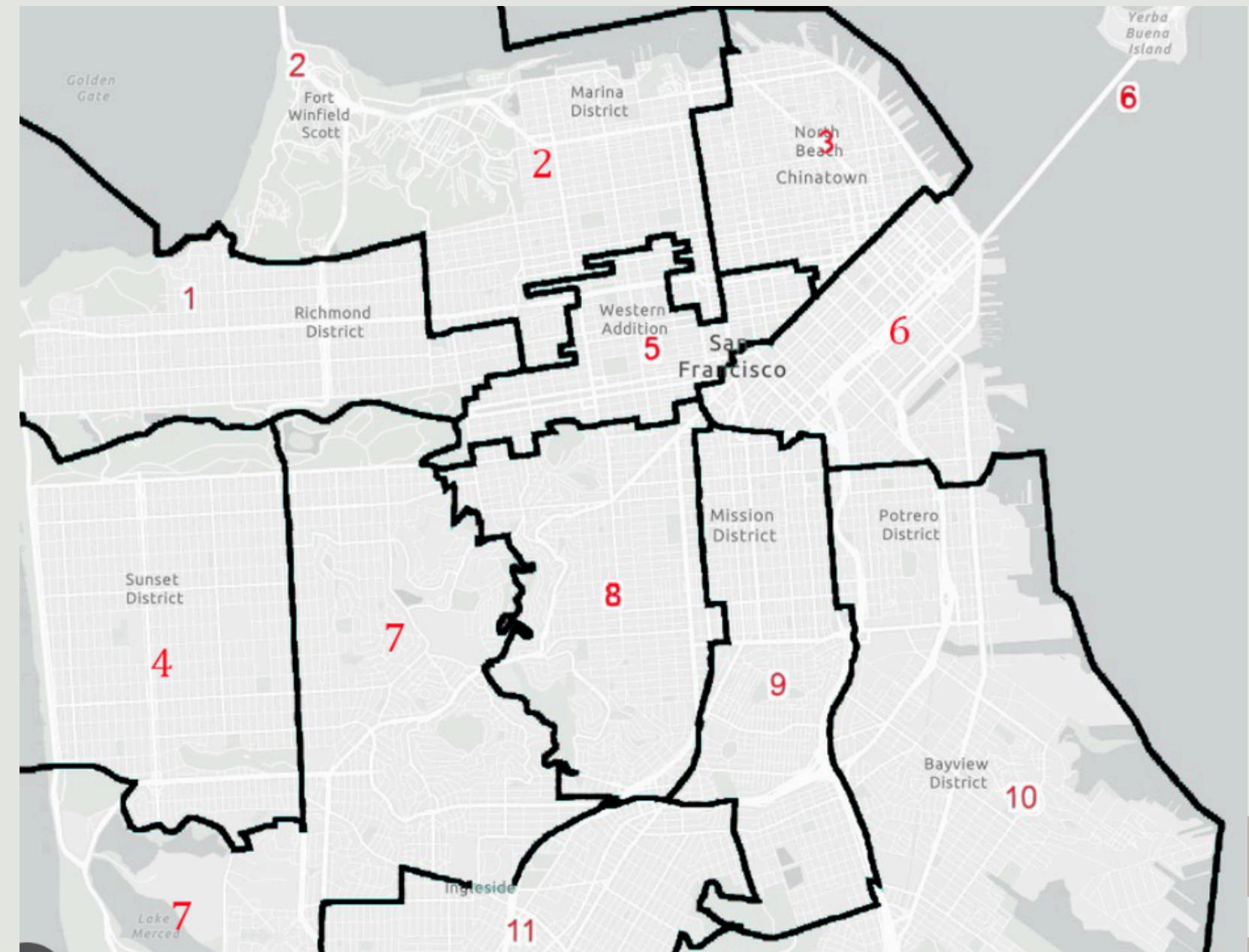
- A Python script collects data on registered businesses from the [San Francisco Open Data API](#) and saved in a CSV format for further analysis.

Web API

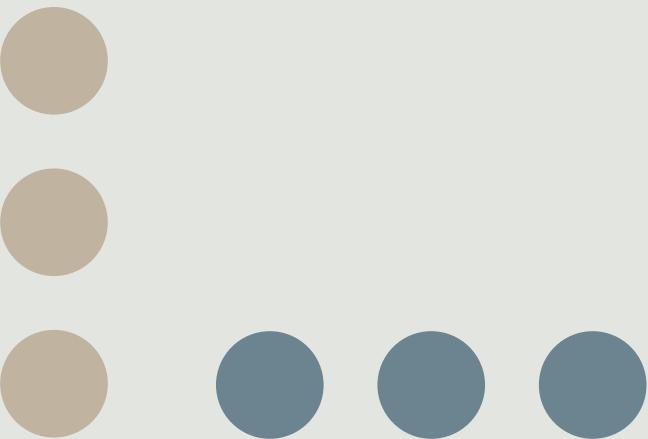
- Endpoint:
<https://data.sfgov.org/resource/g8m3-pdis.json>
- Limit per request: 1,000 records (as per API specifications)

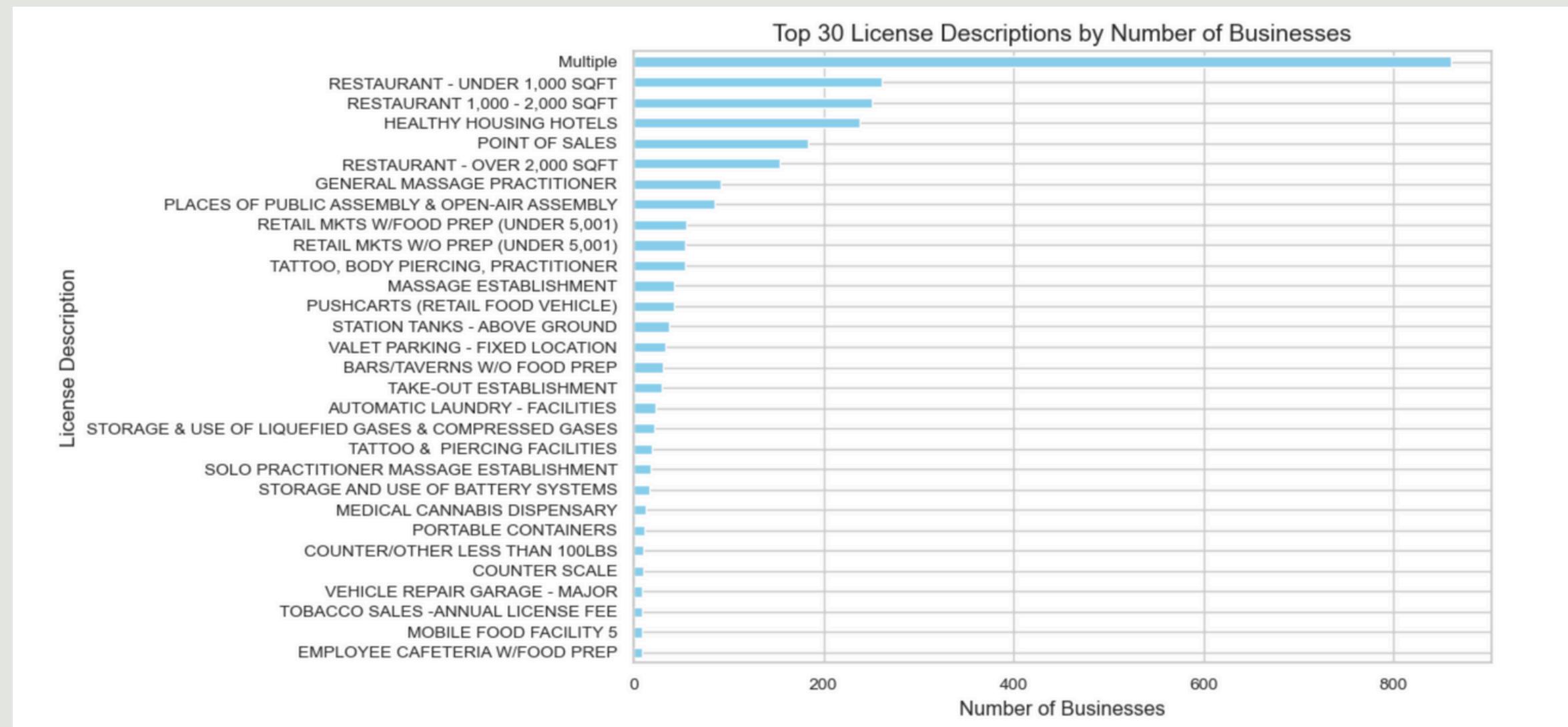
DATA CLEANING AND PREPROCESSING

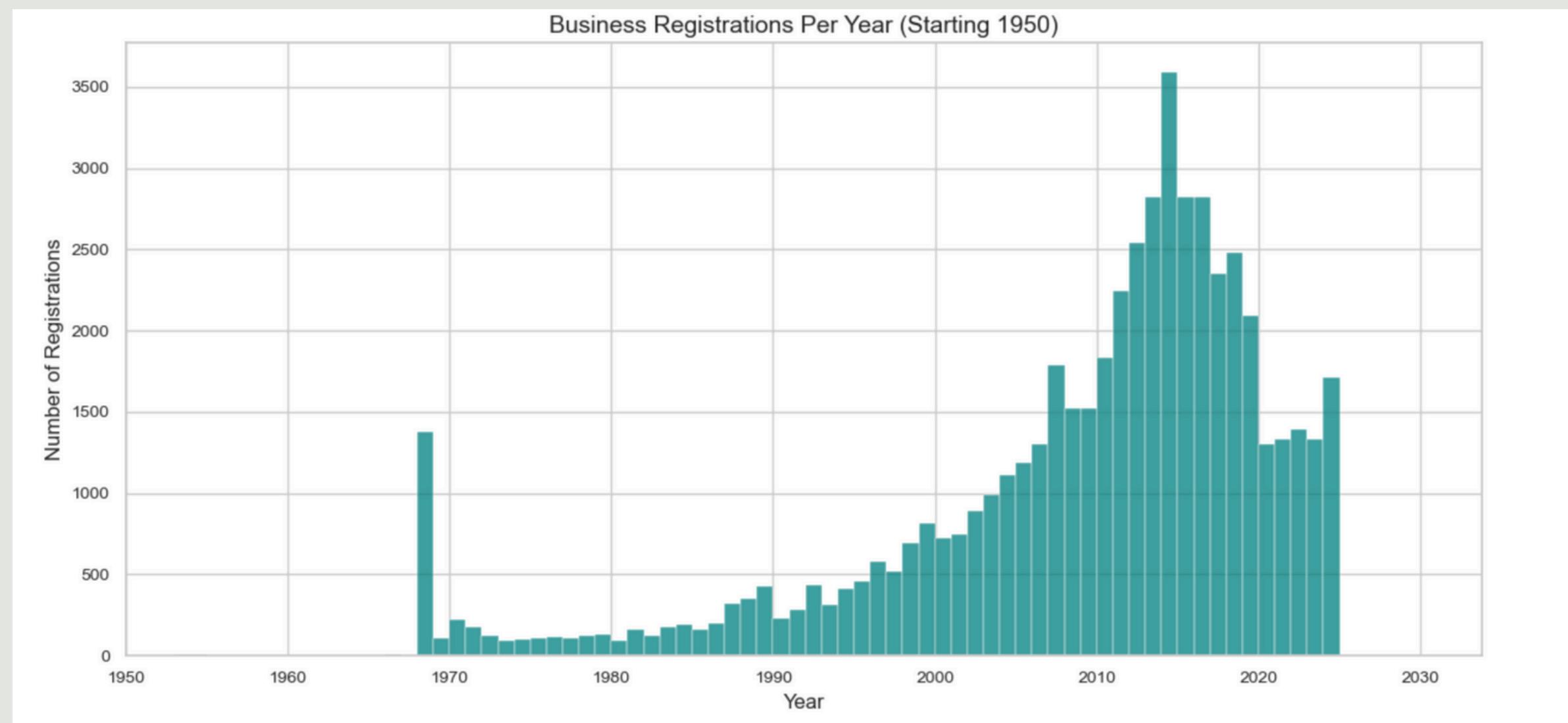
- Standardizing City Names
- Correcting Misspellings of "San Francisco"
- Renaming Columns
- Filtering Data for San Francisco District 3
- Converting Dates and Calculating Business Age
- Duration variable: `business_age` – represents how long a business has been operating.
- Event indicator: `administratively_closed` – indicates whether a business was closed (1) or is still operational (0).

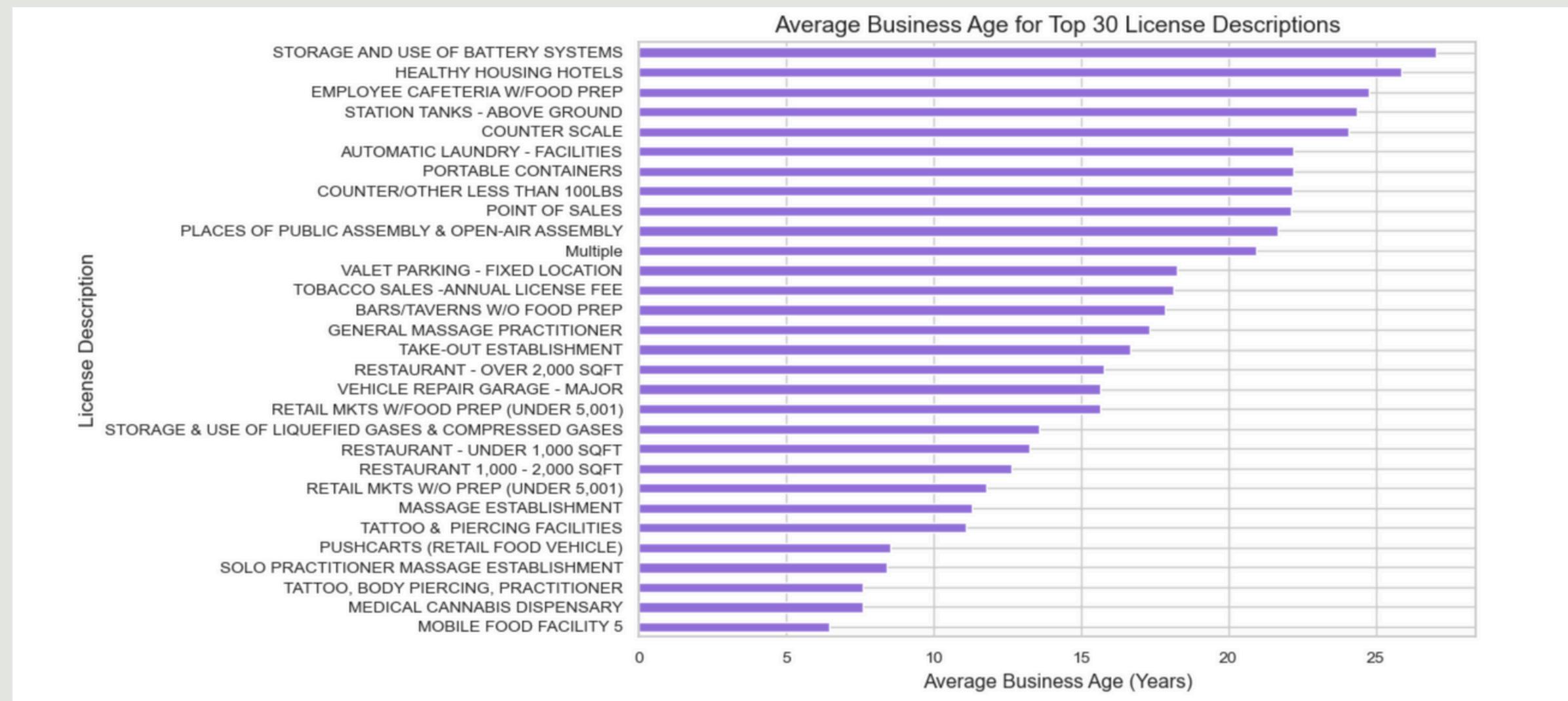


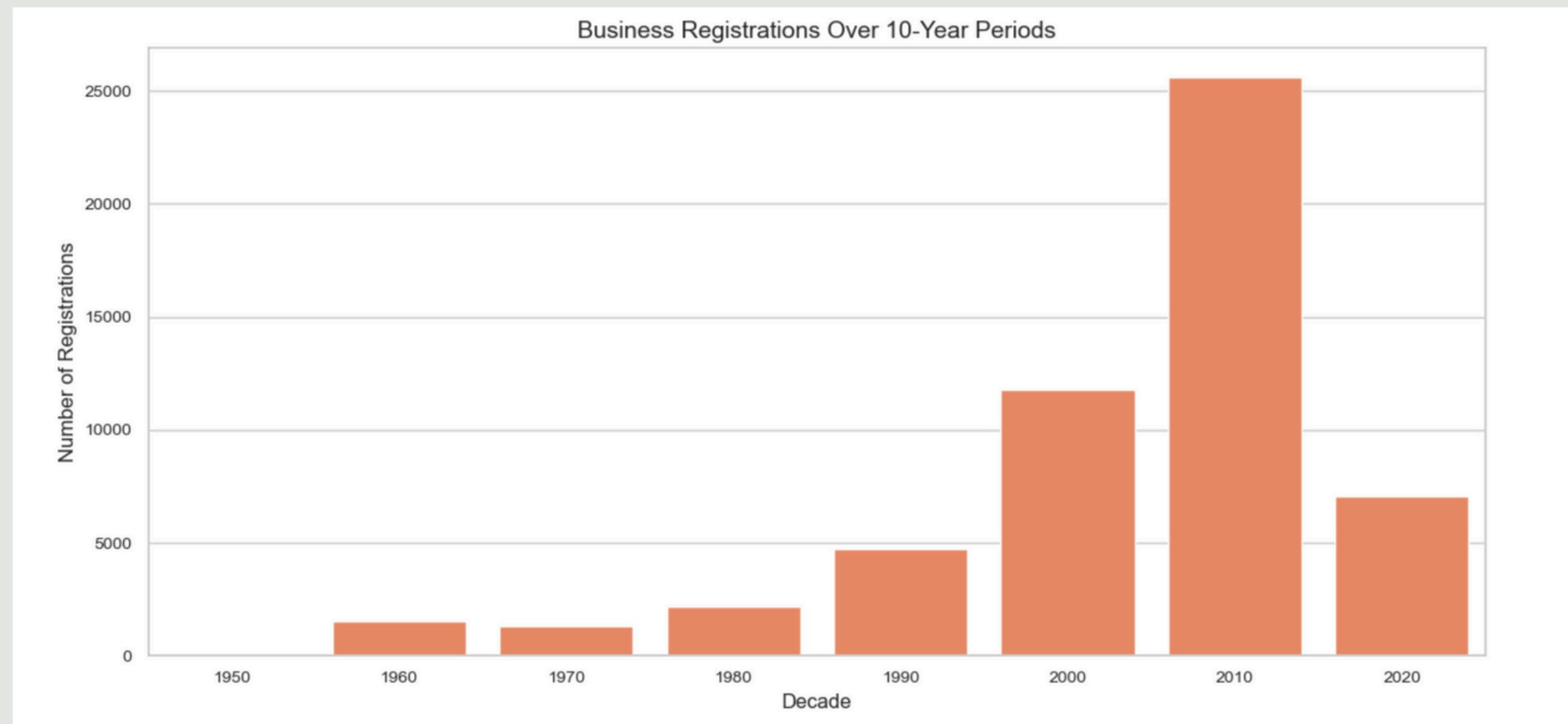
EXPLORATORY DATA ANALYSIS (EDA)

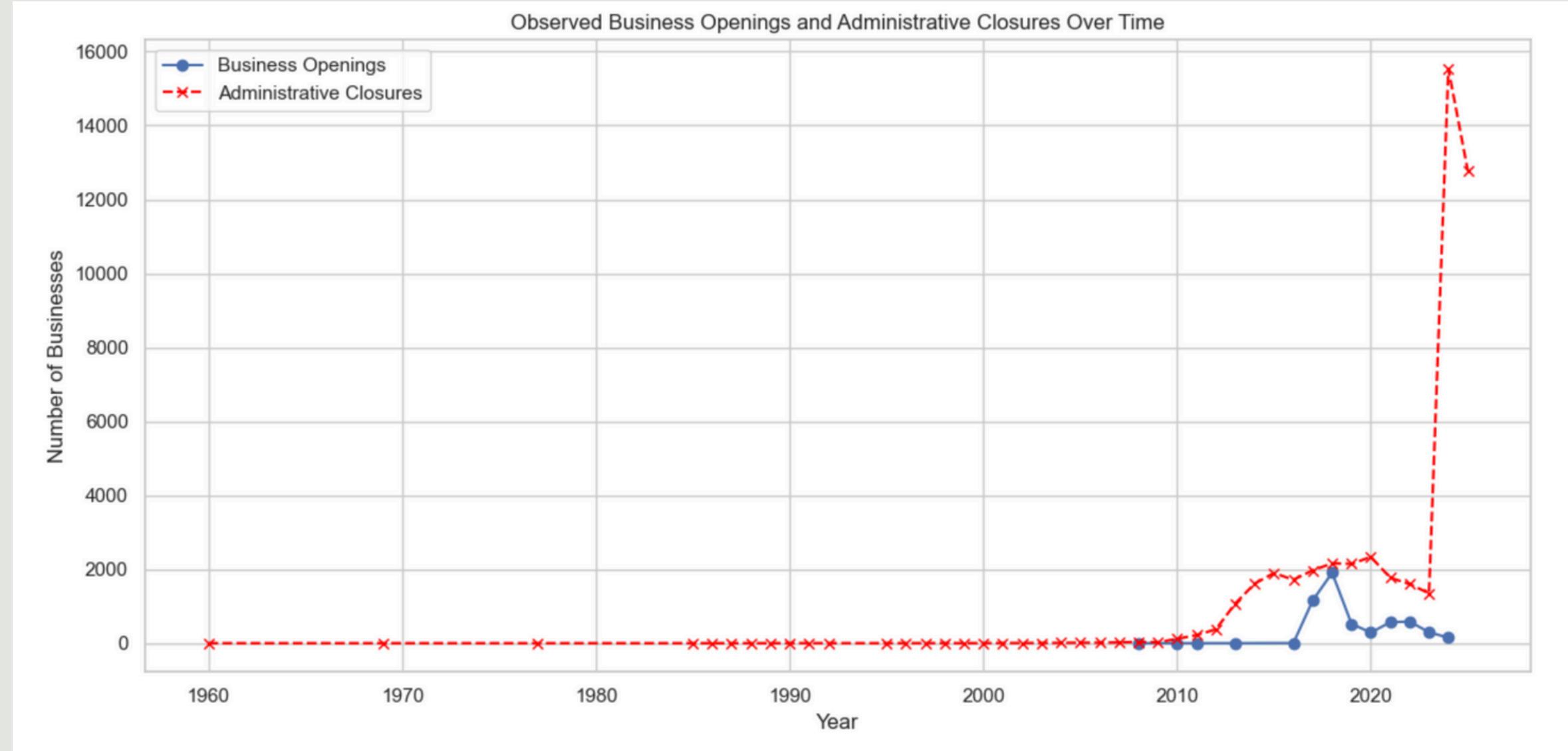












MODELING

- Model: Cox Proportional Hazards
- Features: Business age, zip, tax fields, NAIC codes
- Evaluation: Concordance index
- Findings:
- Newer businesses = higher risk
- Tax codes & zip codes matter

FEATURE IMPORTANCE AND MODEL EVALUATION

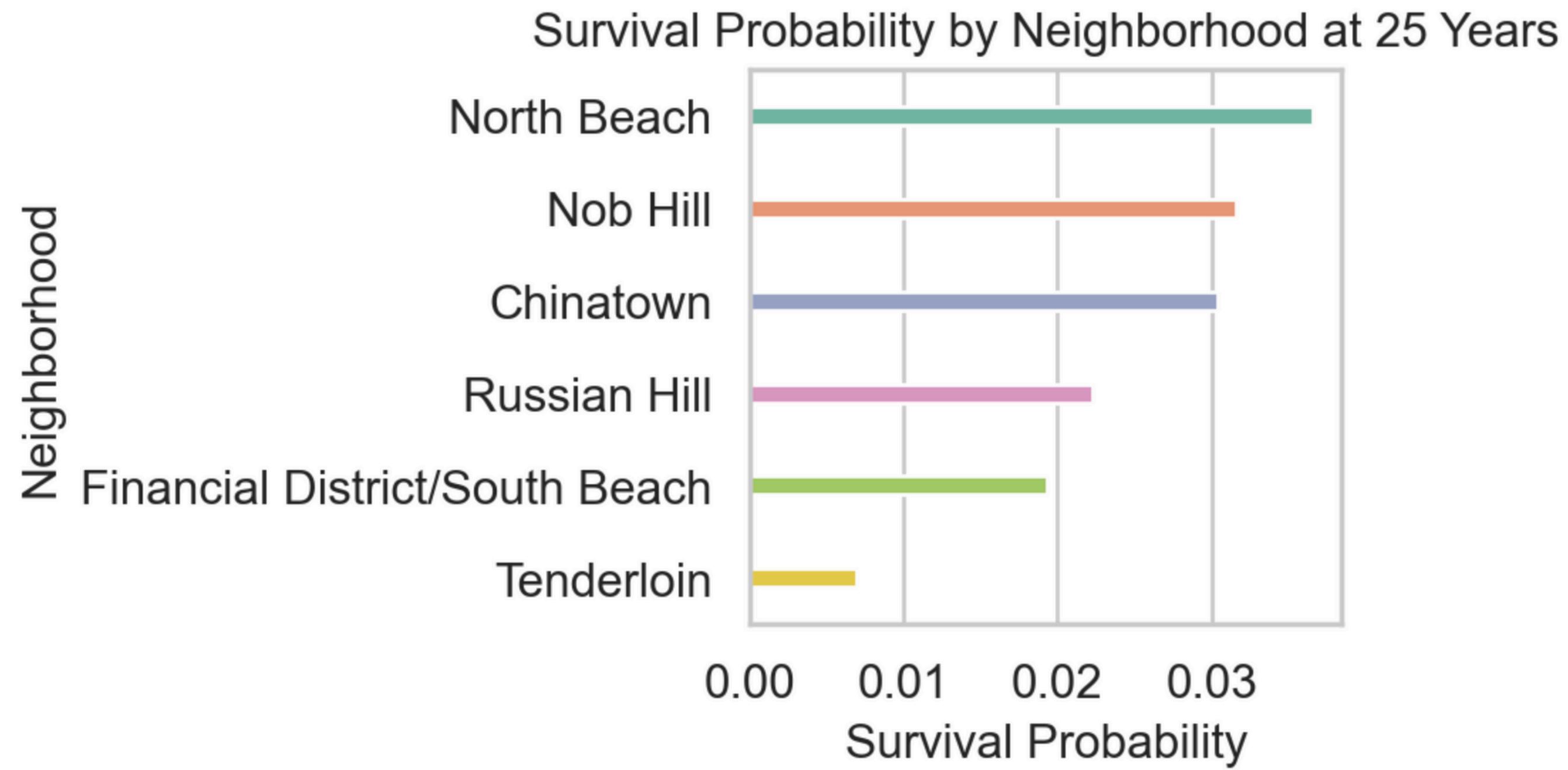
Principal Component Analysis:

- Applied PCA Principal Component Analysis to find the feature importance to improve model performance
- Model performance: Concordance = 0.74

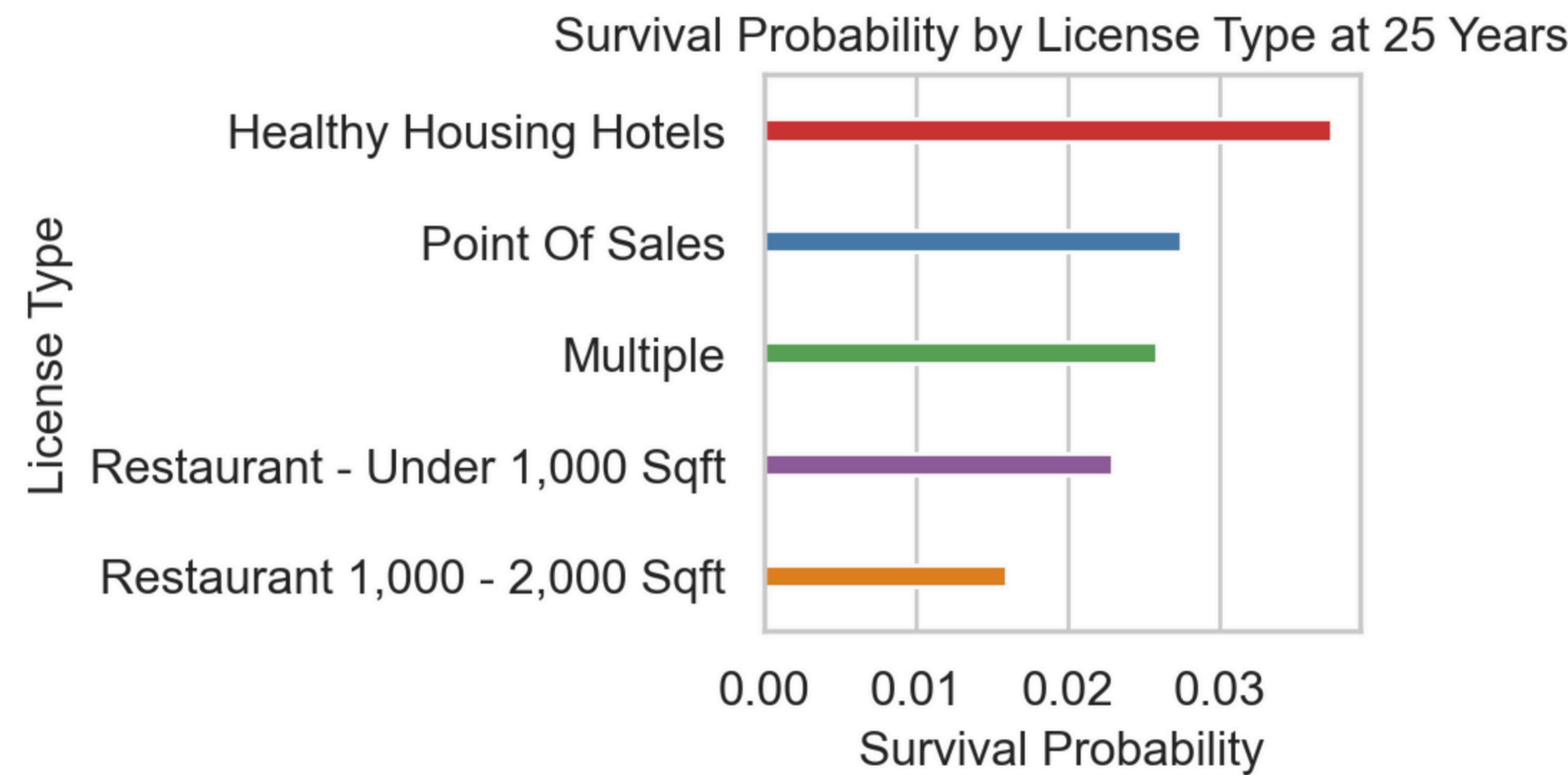
Top Contributing Features:

- Increase in Transient Occupancy Tax/Parking Tax cause higher business Closure risk
- Neighborhood & Zip has a moderate influence in business closure

🔍 Survival Probability at 25 Years by Neighborhood



 Survival Probability at 25 Years by License Type



SURVIVAL TRENDS

- 5-Year Survival (Avg): 11.03%
- 10-Year Survival (Avg): 7.23%
- Best Neighborhood: Russian Hill (14.81%)
- Best License Type: Restaurant (12.75%)

RECOMMENDATIONS

Neighborhood-Based

- Identify underperforming neighborhoods using survival rates
- Provide targeted grants to support vulnerable areas
- Invest in infrastructure improvements to boost business activity
- Offer mentorship and advisory programs for new or struggling businesses

License-Type-Based

- Analyze top-performing license types to identify best practices
- Use successful categories as models for at-risk business types
- Design targeted support programs for struggling license types
- Consider policy changes and financial incentives to improve survival outcomes

Streamlit

Demo



References

Number Analytics – Kaplan-Meier Survival in 5 Steps

🔗 <https://www.numberanalytics.com/blog/kaplan-meier-survival-5-steps>

Lifelines Documentation – Survival Regression

🔗 <https://lifelines.readthedocs.io/en/latest/Survival%20Regression.html>

SF Open Data Portal – DataSF

🔗 <https://datasf.org/opendata/>

Questions?

Thank you