

1 Automated data extraction from historical city directories:  
2 the rise and fall of mid-century gas stations in Providence, RI

3

4 Samuel Bell,<sup>1\*</sup> Thomas Marlow,<sup>2</sup> Kai Wombacher,<sup>1</sup> Anina Hitt,<sup>3</sup>

5 Neev Parikh,<sup>3</sup> Andras Zsom,<sup>1</sup> Scott Frickel<sup>2</sup>

6

7 <sup>1</sup>Advanced Research Computing, Center for Computation and Visualization, Brown  
8 University, Providence, Rhode Island, United States of America

9

10 <sup>2</sup>Institute at Brown for Environment and Society, Brown University, Providence, Rhode  
11 Island, United States of America

12

13 <sup>3</sup>Brown University, Providence, Rhode Island, United States of America

14

15

16 \*Corresponding author

17 Email: [samuel\\_bell@brown.edu](mailto:samuel_bell@brown.edu)

18

## 19 **Abstract:**

20

21 The location of defunct environmentally hazardous businesses like gas stations has  
22 many implications for modern American cities. To track down these locations, we  
23 present the *directoreadr* code ([github.com/brown-ccv/directoreadr](https://github.com/brown-ccv/directoreadr)). Using scans of Polk  
24 city directories from Providence, RI, *directoreadr* extracts and parses business location  
25 data with a high degree of accuracy. The image processing pipeline ran without any  
26 human input for 94.4% of the pages we examined. For the remaining 5.6%, we  
27 processed them with some human input. Through hand-checking a sample of three  
28 years, we estimate that ~94.6% of historical gas stations are correctly identified and  
29 located, with historical street changes and non-standard address formats being the  
30 main drivers of errors. As an example use, we look at gas stations, finding that gas  
31 stations were most common early in the study period in 1936, beginning a sharp and  
32 steady decline around 1950. We are making the dataset produced by *directoreadr*  
33 publicly available. We hope it will be used to explore a range of important questions  
34 about socioeconomic patterns in Providence and cities like it during the transformations  
35 of the mid-1900s.

36

## 37 **1. Background**

38

39 Until the passage of the Resource Conservation and Recovery Act (RCRA) of 1976,  
40 waste produced during commercial and industrial activities in the United States was  
41 largely unregulated (1). It took another decade still until programs such as the  
42 Environmental Protection Agency's (EPA) Toxic Release Inventory (2), were  
43 established for keeping track of emissions from the largest and most hazardous  
44 facilities. These regulatory dynamics, combined with businesses tendencies to  
45 constantly churn in and out of operation over time, has created an urban environment  
46 covered with the relic sites and toxic legacies of past economic activity (3). This is a  
47 serious concern for both community members worried about their health (4), but also  
48 regulators and environmental professionals interested in locating and remediating  
49 contaminated sites. Thus, developing tools for the collection of data on the historical  
50 locations of businesses prone to producing waste represents an important contribution  
51 and is the goal of the current work.

52

53 Previous work focused on developing a software pipeline for processing historical  
54 directories specific to industrial manufacturing (5). The result was a software named  
55 *georeg* (<https://github.com/brown-ccv/georeg>), which was able to process digitized  
56 industrial directories to produce a near comprehensive dataset of industrial site  
57 locations and activities in Rhode Island for the years 1953-2012. This has been used  
58 productively for a range of scientific and community activities (6). However, while  
59 industrial production is a major source of urban pollution - it represents only a selection  
60 of economic activities that leave behind on-site contaminants. Gas stations are another  
61 such commercial activity of concern. According to the EPA, underground gas and oil

62 storage tanks at these sites are a leading source of groundwater contamination (7).

63 And while the federal government has been monitoring underground storage tanks

64 (USTs) since the mid 1980s through RCRA, older USTs are added to lists only as they

65 are discovered.

66

67 Therefore, in the current paper we develop an approach to collecting the historical

68 location of commercial sites from city directories. Since the 1930s, the Polk Corporation

69 has maintained detailed city directories for most American cities. Compiled annually,

70 these books contain a comprehensive list of area businesses in the yellow pages.

71

72 Because the structure of the data in the city directories was considerably different from

73 the industrial registries, we developed a new code, *directoreadr*, instead of adapting the

74 *georeg* code of (5). We did use the same custom geocoder as (5), although the

75 geocoding processing code was quite different.

76

77 To develop and test *directoreadr*, we have focused on city directories from Providence,

78 Rhode Island. Using the scanned images of these directories, *directoreadr* is able to

79 extract a company's name, address, and business type. The data are then geocoded to

80 provide latitude and longitude. To show an example use of these data to examine

81 environmentally hazardous sites, we focus on gas stations. However, the applications

82 are not limited to tracking environmentally hazardous sites. These data can answer

83 many important research questions across a range of topics that are of interest to many

84 social science and environmental disciplines, from economics to ecology.

85

## 86 2. Data

87

88 For the purposes of this project, we have focused on the yellow pages business  
89 directory within the city directories. We examined 27 city directories from the city of  
90 Providence, RI with dates from 1936 to 1990. Beginning in 1940, the city directories  
91 were produced by the Polk corporation, but the three city directories from before 1940  
92 were produced by the Sampson and Murdock corporation. By extracting this detailed  
93 spatio-temporal business data, we allow for socioenvironmental analysis of changes in  
94 the land use of industrial sites, manufacturing zones, or other potentially hazardous  
95 areas, such as current and former gas station sites.

96

97 Digitization was performed by the Internet Archive's office at the Boston Public Library,  
98 and the physical books were supplied by both the Boston Public Library and the  
99 Providence Public Library. The Internet Archive uses a standardized digitization  
100 process, delivering 300 dpi 8-bit color images with a lossy compression in the wavelet-  
101 based JPEG 2000 format. We convert these files to grayscale and do not use color  
102 information.

103

## 104 3 Methods

105

106 The *directoreadr* pipeline consists of a series of discrete processing steps that convert a  
107 series of page images into a database of businesses and locations. These steps are:  
108 grayscale thresholding, ad removal, margin cropping, column chopping, line chopping,  
109 Optical Character Recognition (OCR), header identification, entry concatenation, text  
110 cleaning, address parsing, street matching, and geocoding.

### 111 3.1 Image preprocessing

#### 112 3.1.1 Grayscale thresholding

113 The original color images are read into *directoreadr* as 8-bit grayscale images with an  
114 integer pixel value ranging from 0 to 255, and the first step of the pipeline is to convert  
115 these images to a binary format, where each pixel value is either 0 or 1. This  
116 binarization step enables us to detect connected areas of black pixels, a core  
117 component of many of the computer vision algorithms we use. To do this, we use a  
118 fixed threshold across the entire page. Initially, *directoreadr* will attempt to estimate the  
119 threshold from the distribution of pixels at different pixel values. Incorrect grayscale  
120 thresholds are one of the largest sources of error in *directoreadr*, and its improvement  
121 would provide more accurate results. For this reason, we have chosen to allow a  
122 manual override of the grayscale threshold, and we recommend adjusting the grayscale  
123 threshold of pages with high error rates by hand.

### 124 3.1.2 Ad removal

125 After producing the binary images, we remove the advertisements along the border of  
126 the pages, as well as lines and decorations within and between the columns of text. For  
127 the sake of simplicity, we refer to all page features to be removed as “ads.” To identify  
128 and separate out the ads from the text, we leverage two different geometric  
129 characteristics: First, the ads tend to be outlined by simple shapes like direct  
130 rectangles. Secondly, the ads tend to be much larger in extent than the characters of  
131 the text. Figure 1 shows an example of ad removal.

132

133

134 Figure 1: An example city directory page showing the ad removal process. The first  
135 panel shows the city directory page as a binary image. The second shows the contours  
136 identified as ads by the ad removal algorithm. The final panel shows the image after ad  
137 removal.

138

139 Using the OpenCV *contours* method, we identify regions of connected pixels. For each  
140 pixel contour, we calculate both the perimeter of the contour and the perimeter of the  
141 bounding box, the smallest possible horizontal rectangle circumscribing the contour. In  
142 most cases, ads can be separated from text simply by looking at the perimeter of the  
143 bounding box. However, in a few cases, when the grayscale threshold has been  
144 incorrectly set, many characters of text blur into each other, and the perimeter of the  
145 bounding box can be as large as it is for the smallest ads. To address this, we multiply

146 the perimeter of the bounding box by the ratio between the perimeter of the bounding  
147 box and the perimeter:

148

$$149 \quad p_{adjusted} = \frac{p_{bounding}^2}{p}.$$

150

151 Because text has a more complex shape than the ads, the ratio of bounded perimeter to  
152 contour perimeter is much lower for text than for ads, and it helps separate the text from  
153 the ads. Once we have identified the contours around the ads, we remove any black  
154 pixels within the bounding box around those contours.

155

156 In most cases, the ads around the edge of the page are surrounded by horizontal  
157 rectangles. We can address most cases that are not by identifying where the columns  
158 of ads are and removing all black pixels there. Even then, in a few cases ad removal  
159 fails, and the image has to be cropped by hand.

160

### 161 3.1.3 Margin cropping

162

163 Once the ads have been removed and replaced by whitespace, the columns of text in  
164 the center of the page are still surrounded by whitespace. To focus in on the text, the  
165 next step in the *directoreadr* pipeline is to remove the whitespace. To allow for specks,  
166 lines, and other noise on the page, we set a pixel threshold. Margin cropping is fairly  
167 straightforward and rarely creates problems.



168

## 169 3.2 Image segmentation

### 170 3.2.1 Column chopping

171

172 Each page is set up with columns of text (usually three columns), and in order to  
173 preserve information about text location, we separate the text into the columns. To  
174 identify the column breaks, we sum up the number of black pixels in each vertical line of  
175 pixels in the image. Around the column breaks, there are dips in the number of black  
176 pixels. To identify the location of the dips, we set a pixel threshold and identify the  
177 vertical lines with fewer black pixels than the threshold value. We cluster those vertical  
178 lines using the mean-shift machine learning algorithm. Unlike traditional clustering  
179 algorithms, like k-means, which take a number of clusters as an input, mean-shift  
180 figures out the optimal number of clusters. As the cut point for column separation, we  
181 pick the right-most vertical line in each cluster.

182

183 One of the key features of this algorithm is to err on the side of failure, throwing an error  
184 when the ad removal has performed poorly. The goal of this design is to allow for hand-  
185 chopping when it will meaningfully improve the results, and for all of the failure cases,  
186 we generated the columns through hand-chopping.

187

## 188 3.2.2 Line chopping

189

190 Once we have the columns of text, we then chop the columns into individual lines of  
191 text. To identify the lines of text, we use a similar process to identifying the columns.  
192 We calculate the number of pixels in each horizontal line of pixels in the column. Then,  
193 we cluster the horizontal lines of pixels that fall below the pixel threshold, using mean-  
194 shift to identify the entry breaks. If there are large blocks of entries that don't separate,  
195 we then run the algorithm on them with a higher black pixel threshold. This higher  
196 threshold is typically necessary when the page image is warped or tilted. This process  
197 is highly robust, and it rarely produces errors unless there are more serious problems  
198 with the image.

199

## 200 3.3 Address processing

### 201 3.3.1 OCR

202

203 Entering this part of the pipeline, we have a directory of images where each image  
204 represents a single line from one of the columns on the page. To convert these images  
205 of text to a string of text, we use the Tesseract OCR package developed by the Google  
206 corporation (8). OCR is not perfect, and it does produce some errors, and downstream  
207 text parsing parts of *directoreadr* must account for these errors.

208

### 209 3.3.2 Header determination

210

211 The data in the city directories are grouped under headers that describe the type of  
212 business, and these headers must be identified. Depending on the year, the city  
213 directories identify headers using a number of different characteristics. Headers are  
214 typically indented, and they sometimes contain all caps. Often, headers have asterisks  
215 before them. Depending on the year, *directoreadr* selects from five different header  
216 determination algorithms, most of which center around how many pixels each line is  
217 indented by. Because some columns are tilted, we calculate a relative indentation  
218 compared to nearby lines. Our header detection algorithms relied on indentation and  
219 capitalization as the primary detection features, and we did not build a robust header  
220 algorithm for 1964, the one year in which headers were not indented, and both the  
221 headers and the text were in all caps. As a result, most header identifications for 1964  
222 are incorrect.

### 223 3.3.3 Entry concatenation

224

225 In many cases, entries in the columns of text are too long for one line and continue onto  
226 the next line. In all of these cases, the next line is indented, but not by as much as a  
227 header is. Using the indentation data, we concatenate the multi-line entries into single  
228 strings of text.

229

### 230 3.3.4 Text cleaning

231  
232 Most of the raw entries just contain a business name and an address, but some of them  
233 contain additional information that must be removed, like a telephone number or a floor  
234 or room number. To clean these data, we used a complex series of pattern matching  
235 operations. In some cases in older books, there were multiple addresses for a business  
236 in a single entry, and we split these lists based on the positions of commas and the  
237 word “and.”

238

### 239 3.3.5 Address parsing

240  
241 We start by using a regex to search the entry’s string of text for abbreviations like:St.,  
242 Av., Ct., Dr., Rd., Ave., and Ln. in either upper or lower case. If one of those  
243 abbreviations is detected in the string, the algorithm searches for a group of digits  
244 before the abbreviation. It then classifies the string of text between the number and  
245 abbreviation the address and the text before the address number is classified as the  
246 company name. If the abbreviation is not detected, the algorithm will still try to parse out  
247 the address by searching for the address number and classifying the string of text after  
248 the number as the address, and string proceeding the number as the company name.

249  
250 This parsing algorithm is not perfect – i.e. it requires digits (not spelled out numbers) for  
251 the address. However, it is generalized enough to work well across many different

252 formats because it is built on simple components of the address that are consistently  
253 present.

254

### 255 3.3.6 Street matching

256

257 Because a number of the streets contained OCR errors, we used fuzzy matching to  
258 produce true street names. We developed two lists of streets, a list of current streets  
259 and a list of historical streets. The historical street list was developed through hand  
260 examination of historical maps and is not fully comprehensive. Because we only had a  
261 database of Providence streets, we removed the addresses we could identify as  
262 belonging to another Rhode Island municipality.

263

264 Using the *fuzzywuzzy* package in Python, we created a scoring algorithm to quantify  
265 how close an OCR reading of a street name is to a street in the true street name list.

266 This scoring algorithm is based off the Levenshtein distance ratio:

267

$$268 \quad \text{ratio}(s1,s2) = 100 * \left(1 - \frac{D(s1,s2)}{L(s1) + L(s2)}\right),$$

269

270 where  $s1$  and  $s2$  are the two strings being compared,  $L$  is a function giving the length of  
271 a string and  $D$  is a function giving the Levenshtein distance between two strings. The  
272 Levenshtein distance is the minimum number of edit operations (substitutions,  
273 deletions, or additions) required to convert one string into another. For instance, the

274 Levenshtein distance between “park” and “barks” is 2, one substitution and one  
275 addition. The ratio would be  $100 \cdot (1 - 2/9) = 77.8\%$ .

276

$$277 \quad \text{score} = \frac{\text{ratio}(s1,s2) + \text{ratio}(\text{LongestWord}(s1),\text{LongestWord}(s2))}{2},$$

278

279 where *LongestWord* is a function that gives the longest word of a string. The reason for  
280 adding additional emphasis on the longest word was to emphasize the core street  
281 name. For instance, we wanted “BROADWAY” to match with “BROADWAY ST.”

282

283 For each street in the true streets list, we calculated the matching score between the  
284 OCR result and the street in the true streets list, selecting the true street with the  
285 highest score. To guard against false positives, we removed any matches with a score  
286 below 80.

287

288 To improve the efficiency of the street searching, we only searched for unique OCR  
289 results and built up a dictionary of OCR results and their cached street search results.

290

### 291 3.3.7 Geocoding

292

293 The last component of the *directoreadr* pipeline is geocoding the cleaned and parsed  
294 addresses to obtain the latitude and longitude coordinates of the businesses. After  
295 researching several different geocoding options, from paid services (SmartyStreets) to

296 free APIs (Google Maps), we decided to implement the geocoder built for (5) using  
297 ArcGIS software, with data from Rhode Island's E911 database. This geocoder was  
298 free for our use, given that Brown University had an in-house ArcGIS server, but  
299 because geocoders are proprietary, in our published version of the code, we do not  
300 include the api key necessary to run the geocoder.

301

302 To improve the speed of the geocoding, we ran 50 concurrent searches and searched  
303 only for unique addresses, building up a dictionary of geocoder results to reference in  
304 future runs of the program. Because many addresses were repeated across many  
305 years, this drastically sped up the process.

306

307 The geocoder only contained data on current street layouts. Providence, however, like  
308 many American cities, has seen considerable change in its street pattern over the  
309 course of the study period. Many streets have been wholly or partially demolished, and  
310 others have been renumbered. To address this problem, we utilized the geocoder  
311 confidence score, and we removed any addresses with a confidence score under a  
312 perfect score of 100 from our final results. In a case where the entire street was no  
313 longer remaining, the geocoder returned an error. To address the most common of  
314 these addresses, we allowed for hardcoding of hand-identified historical geocodes,  
315 entering hard-coded locations for four large buildings with many businesses at those  
316 addresses.

## 317 4. Results and discussion

318

319 The image processing portion of the pipeline had a success rate of 94.4%. In our  
320 dataset, we ran the algorithm on 2,582 individual pages. For these pages, 144 or 5.6%  
321 required hand-chopping in order to process. We designed the column-chopping  
322 algorithm to deliberately fail when there were likely errors with the ad removal algorithm.  
323 The goal was to require hand-chopping whenever it would meaningfully improve the end  
324 result. Because of the hand-chopping, we were able to pass all of the pages through to  
325 the OCR and text parsing algorithms.

326

327 In the text parsing algorithm, 6.7% of all entries were dropped as not a successfully  
328 identified and matched address. These include both entries that should be dropped and  
329 entries that were dropped because of an error. In 38.2% of these cases, the algorithm  
330 failed to parse an address at all. In 10.3% of these cases, the algorithm parsed an  
331 address but returned an empty string for the street. In 4.2% of these cases, the  
332 algorithm parsed an address but threw an error in street matching. And in 47.2% of  
333 these cases, the algorithm successfully parsed an address and matched a street, but  
334 the confidence score was too low for us to be sure the address was correct. In some of  
335 the address drop cases, the addresses were outside of Providence, sometimes outside  
336 of Rhode Island. Others reflected an idiosyncratic address form the algorithm wasn't  
337 set up to parse. For instance, some addresses were named buildings without an  
338 address (e.g. "Arcade Bldg" or "Industrial Trust Bldg"). Others were street corners  
339 instead of numerical addresses. Of course, many of the address drop cases



340 represented failures of the OCR or failures of the header identification, concatenation,  
341 and entry chopping algorithms. Address drop rates were not strongly correlated with  
342 time (Figure 2).

343

344

345 Figure 2: Address drop and geocoder error rates by year.

346

347 The geocoder algorithm produced errors in 4.7% of cases. Unsurprisingly, these errors  
348 were higher in earlier years when the Providence street pattern was considerably  
349 different (Figure 2). Towards the end of the study period, the percentage of addresses  
350 outside of Providence increased sharply. While we were able to capture most of these,  
351 we were not successful in all cases. Many addresses from a different city were not  
352 recognized as belonging to a different city, and when they were processed, they led to  
353 dropped addresses or geocoder errors.

354

355 These statistics only capture the places where the code generated errors or flags. In  
356 order to fully assess the ultimate accuracy of the code, we hand-examined the error rate  
357 for gas stations in three years: 1936, 1962, and 1990. In 1936, 220 out of 242 gas  
358 stations were correctly identified, for an accuracy rate of 90.9%. Of the 22 missing gas  
359 stations, eleven were missing due to geocoder errors cause by historical street  
360 changes, seven were missing because of non-standard address formats that  
361 *directoreadr* could not parse correctly, one had the wrong address read, and only three  
362 were entirely missing. In 1962, 219 out of 224 gas stations were correctly identified, for

363 an accuracy rate of 97.8%. Of the five missing gas stations, one was dropped, one had  
364 a geocoder error, and three had their addresses read incorrectly. In 1990, 71 out of 73  
365 gas stations were correctly identified, for an accuracy rate or 97.3%. Both of the  
366 missing gas stations were dropped. These statistics do not include errors in correctly  
367 reading and parsing the business names.

368

369 The total accuracy rate from the gas stations in the three years we examined by hand  
370 was 94.6%. (This number is somewhat skewed because 1936 had the most gas  
371 stations. The average of the three accuracy rates was 95.3%.)

372

373 We are making the data available for download from the Brown Digital Repository  
374 (<https://doi.org/10.26300/typ4-nj27>), and we are making the *directoreadr* code available  
375 at [github.com/brown-ccv/directoreadr](https://github.com/brown-ccv/directoreadr).

376

## 377 4.1 Efficiency

378

379 Once the geocoding, street matching, and OCR results have been cached, the parsing  
380 algorithm runs in roughly 20s per book on a standard laptop, enabling faster debugging  
381 and development. With no cached results, the full *directoreadr* pipeline still runs in  
382 under 30 minutes per book. Before our efficiency improvements, *directoreadr* would  
383 take many hours to process a single book.

384

## 385 4.3 Example hazardous site: gas stations

386

387 Because of their environmental importance, we selected gas stations as an example  
388 hazardous site type. Gas stations in Providence typically developed along main roads,  
389 avoiding wealthier neighborhoods like Providence's East Side (Figure 3).

390

391 Figure 3: Gas stations mapped in Providence from 1940 to 1990

392

393 Figure 4: Total number of gas stations recorded in Providence by director year

394

395 Starting in the 1950s, gas stations began a precipitous decline in Providence (Figures 3  
396 and 4). By 1990, there were only 75 gas stations in the city, a decline of 71% since  
397 1950, when city directories list 257. This drop corresponds with a decline in the city's  
398 population, which dropped by a third between 1950 and 1980, the combined result of  
399 job loss from deindustrialization and displacement of minority residents whose  
400 neighborhoods were cleared for several ambitious "urban renewal" projects (9–13).  
401 These changes were part of broader national trends of suburbanization and economic  
402 decline in the urban core (14–16). Other factors specific to the service station and  
403 automobile industries also may have played a role. Broader regulatory changes likely  
404 also affected gas station counts, with zoning having a particularly important effect (17–  
405 19). Because the rate of geocoder errors was higher in the earlier years, these figures  
406 probably underestimate the dramatic drop in the number of gas stations. Overall, we  
407 identified 526 unique gas station addresses in Providence over the study period,  
408 compared to just 114 gas station addresses recorded in the Rhode Island Department

409 of Environmental Management Underground Storage Tank (UST) database  
410 (<http://www.dem.ri.gov/programs/wastemanagement/inventories.php>).

411

## 412 5. Conclusions

413

414 We have successfully built a pipeline for the digitization, extraction, and processing of  
415 city directory data. While this approach was developed on directories from Providence,  
416 RI, these directory formats are fairly similar in different cities, and this approach should  
417 be adaptable to cities all across the country. There are many potential uses of these  
418 data, and we have demonstrated mapping of environmentally hazardous historical gas  
419 station sites as an example.

420

## 421 Acknowledgements

422 We acknowledge support from the Institute at Brown for Environment and Society, which funds  
423 a Research Assistantship for T.M., and from the Superfund Research Program of the NIEHS  
424 grant 2P42 ES013660. This work has also benefited from seed grants from the Brown  
425 University Office of the Vice President for Research (grant GR300065) and the Brown Social  
426 Sciences Research Institute. We would like to thank the Providence Public Library and the  
427 Boston Public Library for providing many of the physical directories for scanning, as well as the  
428 Internet Archive for doing the scanning. We would like to specially thank Kate Wells of the  
429 Providence Public Library for her invaluable advice and assistance in arranging the scanning.

430

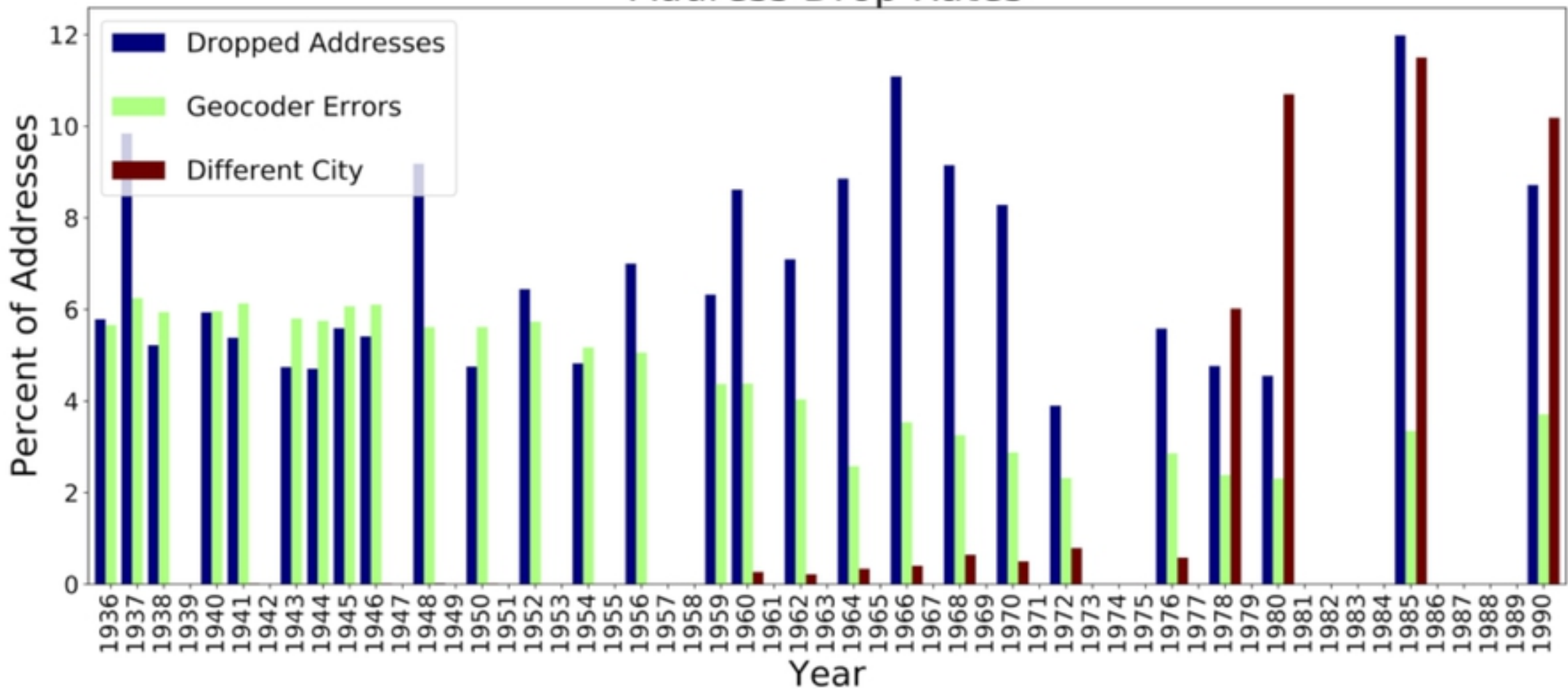
431 Bibliography

- 432 1. Colton CE, Skinner PN. The road to love canal: managing industrial waste before EPA.  
433 Austin: U of Texas P. 1996;
- 434 2. Environmental Protection Agency US. U.S. EPA 1988 Toxics Release Inventory (TRI)  
435 Program. US EPA; 1988.
- 436 3. Frickel S, Elliott JR. Sites unseen: uncovering hidden hazards in american cities. Russell  
437 Sage Foundation; 2018.
- 438 4. Brown P. Toxic exposures: contested illnesses and the environmental health movement.  
439 New York Chichester, West Sussex: Columbia University Press; 2007.
- 440 5. Berenbaum D, Deighan D, Marlow T, Lee A, Frickel S, Howison M. Mining Spatio-  
441 temporal Data on Industrialization from Historical Registries. J ENV INFORM. 2018;
- 442 6. Guelfo JL, Marlow T, Klein DM, Savitz DA, Frickel S, Crimi M, et al. Evaluation and  
443 Management Strategies for Per- and Polyfluoroalkyl Substances (PFASs) in Drinking  
444 Water Aquifers: Perspectives from Impacted U.S. Northeast Communities. Environ Health  
445 Perspect. 2018 Jun 15;126(6):065001.
- 446 7. Environmental Protection Agency US. National Water Quality Inventory: 2000 Report. US  
447 EPA; 2002.
- 448 8. Smith R. An overview of the tesseract OCR engine. Ninth International Conference on  
449 Document Analysis and Recognition (ICDAR 2007) Vol 2. IEEE; 2007. p. 629–33.
- 450 9. Antonucci C. Machine Politics and Urban Renewal in Providence, Rhode Island: The Era  
451 of Mayor Joseph A. Doorley, Jr., 1965-74. 2012;
- 452 10. Goldstein S, Mayer KB. Metropolitanization and population change in Rhode Island.

- 453 Metropolitanization and population change in Rhode Island. 1961;(3).
- 454 11. Goldstein S, Mayer K. Population decline and the social and demographic structure of an  
455 american city. *Am Sociol Rev.* 1964 Feb;29(1):48.
- 456 12. Zimmer BG, Hawley AH. Suburbanization and some of its consequences. *Land Econ.*  
457 1961 Feb;37(1):88.
- 458 13. Zimmer BG. *Rebuilding Cities: The Effects of Displacement and Relocation on Small*  
459 *Business.* Providence, RI: Brown University; 1964.
- 460 14. Jackson KT. *Crabgrass frontier: The suburbanization of the United States.* Oxford  
461 University Press; 1987.
- 462 15. Sugrue TJ. *The Origins of the Urban Crisis: Race and Inequality in Postwar Detroit-*  
463 *Updated Edition.* Princeton University Press; 2014.
- 464 16. Wilson WJ. *When work disappears: The world of the new urban poor.* Vintage; 2011.
- 465 17. Campbell HE, Kim Y, Eckerd A. Local zoning and environmental justice. *Urban Affairs*  
466 *Review.* 2014 Jul;50(4):521–52.
- 467 18. Campbell HE. *Rethinking Environmental Justice in Sustainable Cities: Insights from*  
468 *Agent-Based Modeling.* Routledge; 2015.
- 469 19. Shertzer A, Twinam T, Walsh RP. Zoning and the economic geography of cities. *J Urban*  
470 *Econ.* 2018 May;105:20–39.



# Address Drop Rates





1940

1950

1960

bioRxiv preprint doi: <https://doi.org/10.1101/701136>; this version posted July 12, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



1970

1980

1990

