35 86ae ce81 9b87 f08a 7dda 0d13 5a8c 6674 49e1 bfb0 0153 c8f9 61cd 33ca 6816 66d4 17c8 d222 01a
5e 6eb3 76ab 97e6 9d84 6855 340b 5bec 7e3a 7779 4f6b aa1b 66be 3326 1e47 3dd6 a29d a55e ca4e 4d9
bc ff55 cb0f 19a5 385a 8dcc 1a5b 432e 6296 87d3 a685 ed9a 85f6 f553 916d a530 d79b bed2 9cc0 838
95 cd8f fb65 59af d922 e255 f492 78f0 4677 bf31 2310 c23b 2ed7 bd74 a1e7 78f8 dd19 78ba 4ec4 d96
9f 9623 7cd3 4405 0477 35ca 0801 09d2 153b 8d8a 45ca 5b7c 409b ed93 b697 ff64 45aa 765a 4287 039
c2 6dbc d1da d931 ff6c 0521 1437 628f 521f 1e83 1f24 ec02 8708 49ff 3ce1 45d7 6aa2 e01b 6519 448
56 9259 cb1f 9158 3a5d 9a37 4bb3 f12b 0cdd 9054 f936 8722 b673 7270 10b8 3467 8fe7 6468 a07e c84
ee 1fe9 9f09 9916 decf 963e e4c8 35b4 12c6 770e 2b85 c51d f372 9435 8c34 0101 bd73 2fcd 4725 2ad
d4 a60f 9d4c efce 724f 6e2c a079 799e 99bc f458 c3da 51bd 693f 4e3c d2d1 7fbe 3c4e 5319 a174 e5c
8d 3fd5 a9d8 3af9 62c2 6ff2 3162 db31 03c3 6efb 8b5f b0e4 e82d c0f6 ae67 0b36 9364 dc96 4515 708
a8 2a6d 3922 00d6 e68a bba1 daae c019 c710 eb1a f8e7 254b 7440 0365 db36 bb07 1fba 9dd4 caa9 0ae
aa cdd9 226a 84f1 ea8e a23b 42c4 40d1 d708 d98a 21fd fee0 6f36 8128 61f4 46cd 23aa bd5
df 5f6b b60b e021 d819 21ce d004 6ec8 8    e  81f9 2459 cc1e 9a5e ce53 2eff b8d8 903
44 e577 ddb1 0082 9129 1247 a86c b2a8 02c8 32e1 0030 d133 0bf3 773a 3070 8549 3070 ec94 1de3 d19
3b 9088 7f5b 3e7e f51f 9108 1317 2685 3012 289b e663 73e1 ded3 6833 c95b 4de7 5933 e399 860a 2d1
d0 4dac 8193 b2fb 3165 2b53 d81f 53fa 8ad0 73ca 7481 ea4e 5ffb 85b0 3173 beca 81a0 0915 1d4b 9a5
55 70e9 f1c7 d4db b80c 2e02 6c8a 8823 6da6 5909 2e12 3e30 64a7 a5de 8156 185a 71c1 b6f7 7ce6 41a
4a a994 7307 6fc0 d6e0 f7af d9c6 a7db 4b39 537e 51a6 49ed cd3a f764 8c77 8f00 e737 53c9 f879 4f7
e1 9bae d690 e4d0 6964 b8a7 0842 18dc aacc 3861 df2f f98b 4d68 3be2 4224 ca26 b4fe f421 6834 c90
08 2ab0 190f 3d83 4055 0b1a aeae ab75 7905 0941 eef9 515e 1313 a3c5 69cf 09f0 90e3 886c 5d0d 309
d9 15fd d864 d27e c709 ffa2 77a5 57de 1c6a 0a37 8ee5 805a 0666 0939 0316 de91 ab1a aa23 61d3 086
34 57f8 9234 ce84 c364 5f22 a409 d7a8 eaff 7685 1b06 546e 8d43 d519 4305 df30 4b05 5da8 bfb0 c3c
b7 6eb1 f319 a6a4 ec78 77d3 e8df c9c3 7764 e393 ff94 45fa 5b8d 5296 243a abf6 54e4 64ce 5457 2b8
10 1698 331b 9826 0e38 9c85 1dcb 3fa1 81f4 a7f1 25e7 5b1f 19e6 d7c7 e55c ab20 eeca 75b0 db12 9d2
bb cc04 0b05 825f 5e5a 55f7 06da 6924 ebf0 e676 319d 760f 87de a7cf 67ad 9b36 3f5f 1cea a622 625
69 f4b6 e069 ad69 af12 6550 6efd a4e5 c9f5 9512 1a31 83d7 c600 aadb 6d55 55c3 d6d0 310c 9f8d f3c
b2 69c4 fd63 4476 4e74 b5a1 9c73 1ccd 5827 24ec dde9 fff2 05d5 2abf d141 85eb 387e 2046 5843 4e4
8d f1f3 7e92 c6fe af0f 1df1 5570 477c 52f9 d067 31a1 a8e5 b1fb fe9a 813d 4122 04a6 52c3 e238 7de
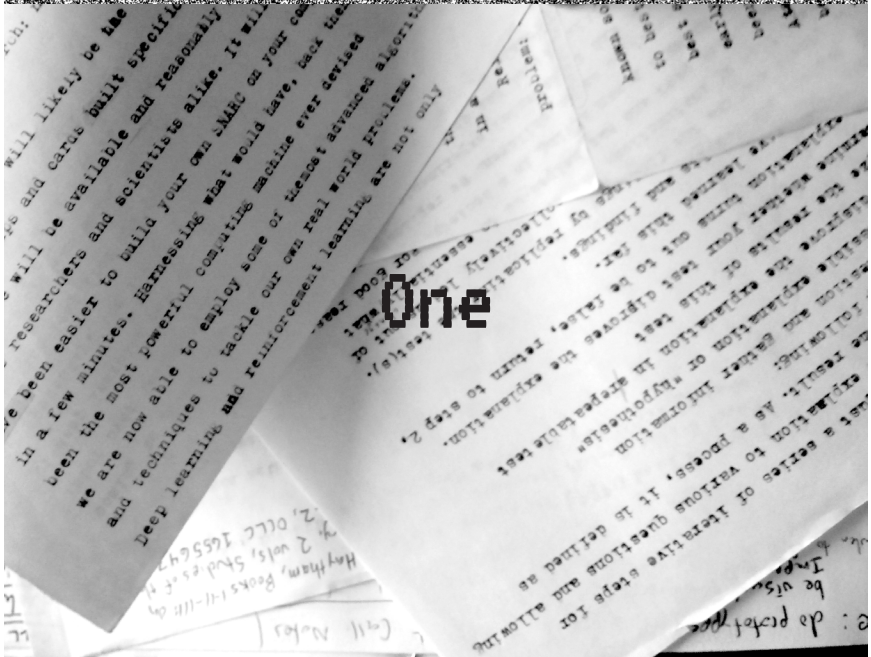f8 b3bf ed88 e04a 9caa 65a2 f837 bfee df2d b027 bd25 a247 91d2 3db8 f250 8953 6304 463d 0b6e ff6

# Issue

# One

# Issue One: Table of Contents

:an handle inexact or "fuz

preprocessing[8] needed

short, LSH operates by t

### Attention!
This is a glossary term! When you see a number in brackets it means the preceeding word has a glossary entry at the end of this zine.

# An Interview with Chase Davis: Journalism + Machine Learning Pioneer

Chase Davis is a journalist based in his hometown, Minneapolis, where he's the Senior Digital Editor at the Star Tribune. He got his start working at The Houston Chronicle, eventually co-founding a media technology consultancy, Hot Type Consulting, which brought him to help launch The Texas Tribune. He's probably best known for leading the Interactive Desk at The New York Times and for being a major presence at NICAR conferences. He's spoken on topics including advanced computational techniques, like machine learning (ML), and innovation at small, local journalism organizations.

We talked to Chase for this first issue not only because he was way ahead of the curve on using ML in the newsroom but also because he has a wide range of practical experience in both the trenches of journalism and in the tech industry. Additionally, he was at The New York Times at a pivotal moment during its storied and successful transition to digital. We reached him via email for this interview. Our questions are in bold.

**Just to get started, can you give us some background on how you got into journalism and where your interest in it comes from?**

I followed a pretty typical path into journalism: I was the editor of my high school paper, went to journalism school at Mizzou, worked at student papers and did a bunch of internships, and kind of moved on from there. As for where the interest came from, I think I just liked the idea of having a career where constant learning was part of the job description.

**As far as I understand, you don't have a formal computer science degree/education. How did you pick up your data and programming skills?**

Not only do I not have a formal computer science degree — I basically failed the only computer science classes I ever took. I started learning to code for fun in middle school, and I kept it up by pursuing personal projects. In college I figured out that using data, basic programming, web scraping, etc., could help me get scoops as a reporter. And then, almost by accident, I discovered IRE and NICAR, which gave me an outlet for refining those skills and applying them more broadly to journalism.

**Your ML presentation at NICAR 2013 with Jeff Larson was really far ahead of its time. Can you explain how you got into ML and applying it to journalism?**

I started getting interested in machine learning pretty soon after I graduated college in 2006. It just seemed like an interesting new frontier with big potential in the world of data journalism. So I started studying it in my spare time, checking out books from libraries, trying to make sense of them, and making little prototypes. At some point I ended up on a couple working groups, notably one put together by Brant Houston at the University of Illinois, between the journalism school and the National Center for Supercomputing Applications there. Those really convinced me there was some potential there.

For a while, machine learning in journalism felt like a solution in search of a problem. But there have been a few problems in recent years where knowing some machine learning has really been helpful. Having some background in the subject has made it easier to step in, see some of those opportunities, and put together solutions that make sense in the context of news.

**A lot of focus is on where ML should be used in journalism. Are there any areas in journalism that you feel ML shouldn't be used?**

Most of them. At this point I think the problem space for supervised learning, in particular, is relatively small and well defined: namely a handful of data cleaning, parsing and classification problems that defy rules-based approaches. Unsupervised learning is more of a question mark. But for most problems in journalism at this point, machine learning is both overkill and even a little dangerous — especially if you don't understand the techniques well enough to adjust for some of their weaknesses.

**Are there any media organizations that you think are doing the best in using these new technologies?**

Beyond the WaPo, NYT, ProPublica, etc., I think the Data Desk at the LA Times has done some incredible stuff.

**Do you see understanding statistics and computer science as becoming a requirement for future journalists, if it's not already?**

We're getting to the point where a basic understanding of things like probability and technology should be a requirement for any professional human, journalist or not. But do I think every journalist needs to be able to code? No. Should they all be able to run a regression analysis? No.

That said, I do think we need to stop with the willful ignorance: "I'm a jour-

nalist, I'm bad at math!" or "I don't do computers!" Not only are assertions like that self-defeating, they sound increasingly absurd as technology's role in society continues to grow.

Plus, our job is basically to learn things for a living. We don't just get to throw up our hands when the learning gets hard.

**It's not ML, but I like to ask people who know about advanced computational techniques about this since journalism is still dominated by it: what do you think about the prevalence of Excel in the newsroom?**

Excel has been the gateway drug for literally thousands of journalists to begin exploring data analysis in the newsroom. Without the right people learning Excel at the right times — and organizations like IRE/NICAR being there to teach them — I'm almost horrified to think of all incredible stories that never would have been told.

**Any advice for people looking to get into journalism in 2019 who can't go to college?**

Do good work and don't be a jerk.

Seriously. In journalism more than most professions, it's all about the work. If you have a demonstrated history of doing good stuff, you're a humble, curious and decent human being, and you plug into the communities of people who do this kind of thing on a daily basis, there's a place for you in this industry.

# Dissecting a Machine Learning Powered Investigation

Uncovering local property tax evasion using machine learning and statistical modeling. An investigative recipe.

By Brandon Roberts

If there's one universal investigative template I've come across in my journalism career, it's this: take a list of names or organizations, find those names in another dataset[1] and identify the most suspicious ones for further investigation. I call this a "data in, data out" investigation. While I've personally only applied this pattern to property tax related stories, it could also be extended to areas like campaign finance and background research.

Before I get into the technical dirt, I want to talk about why data in, data out investigations are ideal. Going into a dataset essentially blind and letting the data provide leads helps keep things less biased. Smaller newsrooms can benefit, too, because they don't need whistleblowers telling them exactly who or what to look for. To me, the best part is these kinds of investigations lend themselves well to advanced computer assisted methods—eliminating a lot of the tedium that goes into investigative work.

I'm going to dissect a data in, data out investigation and the machine learning (ML)[10] tools that went into it.

## The Investigation: An Overview

The City of Austin, like many cities across the US, requires people who list their homes on sites like HomeAway and Airbnb to register and pay taxes on each stay, just like hotels do. These types of rentals are called "short-term rentals" (STRs). Properties used as STRs generally aren't eligible for most tax exemptions. In Texas, home owners can apply for a so-called "homestead exemption". These add up to significant tax savings. Because of this, there are lots of rules around them. For example, only one property

per person or couple can be exempt. It needs to be their primary residence and it can't be used commercially.

We wanted to see whether the appropriate rental taxes were being paid and if the appraisal district was letting STR operators apply for and receive tax exemptions—a clear violation of the law.

A list of STRs registered with the city was the primary dataset. Secondarily, we got datasets of appraisal district properties and electric utility subscribers. We took the STR list and used a ML algorithm to look for matches in the other two datasets. This identified STR operators (individuals or couples) who owned multiple properties and had applied for exemptions.

Next, we checked each joint STR/property record for exemptions and ranked them by suspiciousness. This ranking was accomplished using ML. From there, we filed a series of public information requests to verify our findings.

Joining two datasets based on one of their columns is a basic strategy common to many investigations. Relational databases were created to do exactly this. But in our case, we had low quality data and a lot of it: roughly a million records, combined. Names were spelled inconsistently and addresses were often abbreviated or were missing apartment numbers. Standard tools didn't cut it, so we turned to ML.

Linking Records with Machine Learning

The foundation of my approach to this investigation is a ML algorithm called Locality Sensitive Hashing (LSH). This is part of the unsupervised[19] machine learning family of algorithms. Unsupervised learning gets its name because it learns directly from the data itself. Unlike other forms of ML, unsupervised methods don't require people to tediously sort through data and tag individual records. These algorithms are inexpensive to use, can work in a variety of situations and lend themselves to straightforward, practical applications. Unsupervised learning excels in problems like search, uncovering structure, similarity, grouping and segmentation. It's my opinion that unsupervised machine learning holds the most promise for the future of local and nonprofit investigative journalism.

LSH is an algorithmic method for taking large amounts of data and efficiently grouping together records which are similar to each other. It can handle inexact or "fuzzy" matches, greatly reducing the amount of preprocessing[8] needed.

In short, LSH operates by taking data points and assigning them to "buckets" or groups. When using LSH to link records from two datasets, we take

the larger dataset, build an LSH index and then take the smaller dataset and query that index. The result of a query is a list of records which should be similar to the record we're querying.

I'm going to give a technical overview of the specific LSH algorithm I used. Feel free to skip this part if you don't care.

-------- START TECHNICAL DISCUSSION --------

In order to understand LSH, we first need to discuss the k-Nearest Neighbors (k-NN) algorithm. What k-NN does is take one set of records and allow you to find some number of most similar records. That number is the k in k-NN. So a k-NN algorithm that only finds the single most similar record is known as k-NN with k=1 or (sometimes) 1-NN. In its most basic form, this algorithm checks every point against each other and gets a score for how similar they are. We then rank these and take the k-number of closest matches. While this algorithm has its place, it is problematic in our investigation for several reasons:

> 1. I had a dataset of 20,000 properties and another dataset of 800,000 utilities records. Doing the math, 20,000 records times 800,000, tells us 16,000,000,000 checks need to be performed to find groups using plain old k-NN.

> 2. Even if doing that number of checks was feasible, I'd need to come up with a distance[2] threshold—a number that establishes how similar two records need to be in order to be considered the same.

Both of these problems are annoying to solve and can be avoided by using LSH, which is a form of approximate nearest neighbors.

LSH goes about finding similar records using what I mentally categorize as a mathematical cheat code: if two records, A and B, are similar to a third record, then the two are likely similar to each other. It's a massive simplification of how LSH works, but it serves as a decent mental model. This gets around having to check every record against every record in order to find groups. We only need to check all records against that third thing.

In reality, LSH works by projecting[9] our records onto a randomized hyperplane (basically a random matrix[11] of data), doing a little rounding, and assigning them to a bucket. The smaller the hyperplane is, the higher chance that records will end up in the same bucket, despite dissimilarity.

Generally, this is a lot quicker than k-NN because, for each record, we only need to do a multiplication with the hyperplane. Mathematically speaking, LSH scales linearly with the number of records, as opposed to k-NN which

scales exponentially. In my example above of 20,000 and 800,000 records, I only need to do 820,000 checks (800,000 plus 20,000 checks against our hyperplane). Obviously, this is far less than the 16,000,000,000 checks with plain k-NN.

-------- END TECHNICAL DISCUSSION --------

Grouping records in LSH is done probabilistically, meaning you probably won't find every match and some non-matches might also fall into your results. Fortunately, there are ways to tune LSH in order to avoid this. LSH has two main configuration options, if you will. The most important being the hash[5] size. The larger the hash size, the more similar records must be to end up in the same bucket. Another way to improve accuracy is to use multiple LSH indexes. This increases the chances of records falling into the same bucket if they are indeed similar. I tend to use three LSH indexes with a hash size of either 512 or 1024 bits. (Another strategy for finding a good hash size is to start with the square root of the number of records and go up from there. Wikipedia has a decent page on this.)

The process for grouping records between two datasets using LSH goes like this:

1. Choose a column to match records with and extract that column from each dataset (e.g., we want to match against owner name, so get the "owner name" column from one dataset and "name" on the other).

2. Take only one of our extracted columns for matching and build a LSH index with it (e.g., build a LSH index on the "owner name" column).

3. Take the second, remaining column and query (get bucket assignments from) our LSH index (e.g., for each value in the "name" column, get the most similar names from the "owner name" column). The results of each query forms a "group" of similar records.

4. For each match in each group, grab the original full record. We will use the combined group of records to continue our investigation.

This process results in a list of grouped records. I tend to flatten these groups into wide rows by combining fields from both datasets. I also add some extra columns describing things like the number of matches found in the bucket and some averages from the records. We can use this data to filter out groups that aren't worth investigating further and then rank the remaining records by investigative interest. But instead of doing this last step by hand, I propose using another ML technique: regression.

# Ranking Records by Suspiciousness with Lasso Regression

So, we've successfully identified groups of records and now want to un-cover ones that warrant intensive, manual research. We've also built a few aggregates from each group to work with. In the investigation described above, we had these aggregate variables (or features[4]) to work with:

1. Number of properties owned.
2. Total number of exemptions held.
3. Total monetary amount of combined exemptions.
4. Property type (there were three types of properties, with different rules about how exemptions may legally be applied).
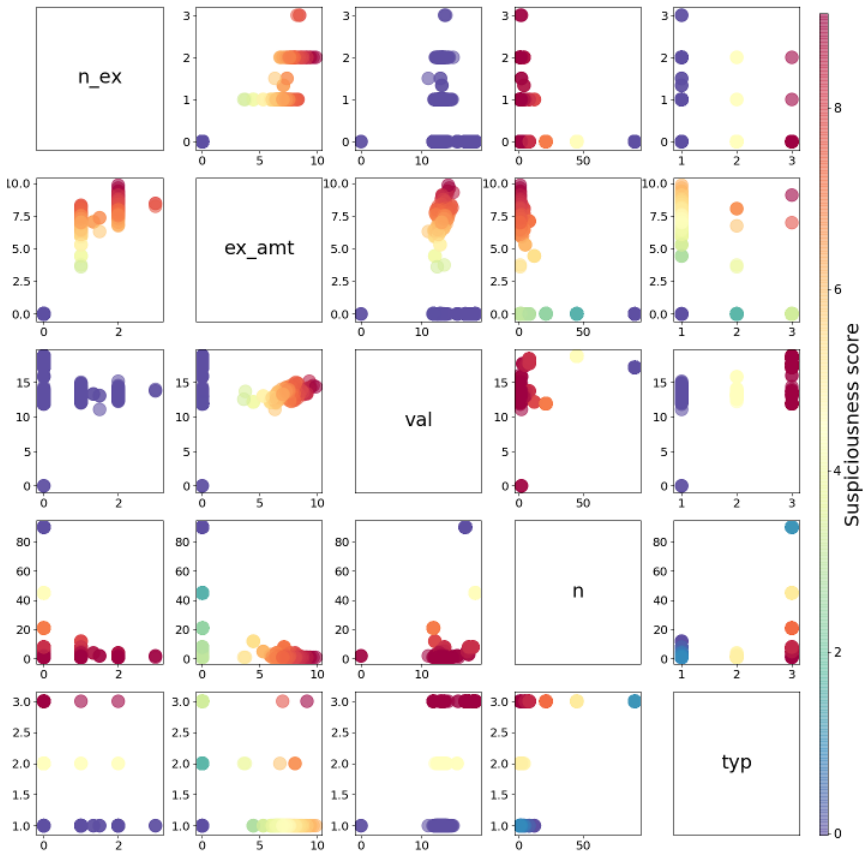5. Average valuation of all properties owned.

Typically, one might consider writing a program to take these values and use them to rank grouped records. We could do this by assigning an impor-tance value, or weight[21], to each variable, combining them and obtaining a score for each group.

This may sound good, but there are still some open questions. How do we combine the variables for each group? How do we compare variables with different ranges? The valuation of a property, for example, is in the hun-dreds of thousands, while the number of exemptions applied to a property is in the single digits. How do we combine everything once we've decided on importance? What do we do if some of the variables are redundant and overpower other, more important, variables? Luckily there is ML algo-rithm that can help with all of these problems.

Regression is a statistical modeling technique in the supervised[17] family of ML algorithms. If you've heard about ML before, it's almost guaranteed to be a supervised learning algorithm. This family of algorithms is useful in cases where you have some task to replicate and examples of it being done. Some problems and algorithms need more examples than others.

A classic example is using regression to predict the price of a home based on the number of rooms, bathrooms, existence of a garage, etc. If we give a regression algorithm enough data about homes along with the price of each (we call this training data[18]), the model can learn to predict price when given information about new homes.

Lasso, short for Least Absolute Shrinkage and Selection Operator, is a specific type of regression that is not only able to make predictions based

*Visual description of a Lasso regression model. In this scatter matrix, each feature is plotted against the other features, forming a scatter plot for each combination. For example, the bottom left plot shows "property type" on the y-axis vs. "number of exemptions" on the x-axis. Points are colored by the suspiciousness scores obtained by the combination of the two features. Red dots are the most suspicious; dark blue/purple dots are the least suspicious.*

on training data, but it can also automatically assign importance values and eliminate useless features from consideration.

In order to use regression to rank our properties by suspiciousness, we need to score some of our records by hand. So I browsed my dataset, found a few seemingly egregious examples, and gave them high scores (I chose zero through ten, with ten being very suspicious). Then I found a few records I wanted to be sure to ignore and gave them a low score. With a few of these hand labeled records, I was ready to build my Lasso model.

There's a major benefit to using regression models: the weights they assign to your input variables are straightforward to understand. Other ML models are basically black boxes. This is so much the case that there's an entire subfield of ML dedicated to building models that can interpret other models. In journalism, we need to quickly check our model's assumptions and biases. Regression gives us a way to do this.

As an example, here's a breakdown of a Lasso model I trained:

```
         Lasso model explanation


Feature                         Importance
_____
Property type           =>  2.3802236455789316
Total exemption amount  =>  0.9168560761627526
Number of properties    => -0.052699999380756465
Number of exemptions    =>  Ignored
Valuation               =>  Ignored


             (Prop type * 2.38)
           - (No. props * 0.05)
       + (Total exempt amt * 0.91)
       ==========================
                Suspiciousness
```

In general, this model thought "property type" was the most important feature. It found records with large numbers of properties owned by a single person worth ignoring. The total exemption amount would bring up a record's suspiciousness score, while number of exemptions and valuation were ignored.

In regression, each of these weights gets multiplied by the corresponding value from each record. When summed, we get a suspiciousness score for each record. Data in, data out.

There were three property types: 1, 2 and 3. Since there were very few type 2 and 3 properties in my dataset, it was important to look into them. Anything rare in a dataset usually is worth looking into. Some of the addresses pointed to large apartment complexes. Unfortunately, our data was lacking apartment or suite numbers, so searches for these properties would bring back every unit in the building. If a bucket contained a lot of properties, it was probably because of a search result like this. Total exemption amount

was also an important variable. In this case, my model figured out that "number of exemptions" and "valuation" were directly related to the "total exemption amount". Because of this, we only needed the total exemption amount and could ignore the other two variables.

Using this, I was able to run my grouped records through the trained Lasso regression model and get back a ranked list of records. The top scoring records being the most likely to be worth manual investigative effort.

# Conclusion

In one investigation, when the algorithms were done spitting out results, we filed public information requests for hundreds of property tax applications. We got a stack of paper four inches thick that my girlfriend and I carefully spread across our living room floor, cataloged and verified. The Austin Bulldog used this analysis to produce a two-part investigative story about the appraisal district's reluctance to vet tax exemption applications. Hundreds of thousands of dollars in back taxes were reclaimed for taxpayers.

In another investigation, we identified properties that were likely running large, illegal online STRs and getting unwarranted tax breaks. The fallout from this one is still ongoing.

Both investigations wouldn't have been possible without algorithmic assistance. But no matter how many fancy computational methods you use during an investigation, you can never escape the need to go back to the primary, physical sources and meticulously check your work.

Machine learning is just a tool, not magic. It won't do our investigation for us and it won't follow the leads it provides. But when used correctly, it can unlock leads from data and help get us to the part we're good at: writing stories that expose wrongdoing and effecting change.

# Rethinking Web Scraping with AutoScrape

By Brandon Roberts

In 2011, living in a city with more dented cars than I'd ever seen in my life, I wanted traffic data for a story I was working on, so I wrote my first real web scraper. I still remember it pretty well: it was written in PHP (barf), ran every five minutes, grabbed the current collisions police were responding to across the city and dumped them into a database. Immediately after writing a story based on the data, I threw the scraper on GitHub and never touched it again (RIP). It got me wondering how many journalists had rewritten and shared scrapers for the exact same sites and never maintained them. I pictured a mountain of useless, time consuming code floating around in cyberspace. Journalists write too many scrapers and we're generally terrible at maintaining them.
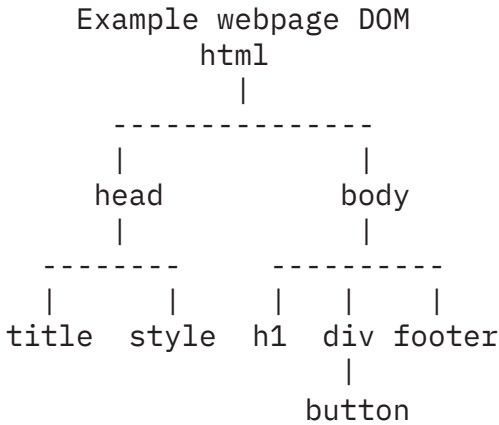
There are plenty of reasons why we don't maintain our scrapers. In 2013 I built a web tool to help journalists do background research on political candidates. By entering a name into a search form, the tool would scrape up to fifteen local government sites containing public information. It didn't take long to figure out maintaining this tool was a full time job. The tiniest change on one of the sites would break the whole thing. Fixing it was a tedious process that went something like this: read the logs to get an idea of which website broke the system, visit the website in my web browser, manually test the scrape process, fix the broken code and re-deploy the web application.

Leap forward to 2018—I was running recurring, scheduled scrapes on open government sites for several media organizations—and I was still building my scrapers the old fashioned way. The fundamental problem with these scrapers was the way they identified web page elements to interact with or extract from: XPath selectors.

To understand XPath, you need to understand that browsers "see" webpages as a hierarchical tree of elements, known as the Document Object Model

(DOM)[3], starting from the topmost element. Here's an example:

```
         Example webpage DOM
                html
                 |
         ---------------
         |             |
        head          body
         |             |
      --------      ----------
      |      |      |   |    |
    title  style   h1  div footer
                         |
                       button
```

XPath is a language for identifying parts of a webpage by tracing the path from the top of the DOM to the specified elements. For example, to identify the button in our example DOM above, we'd start at the top html tag and work down to the button, which gives us the following XPath selector:

```
/html/body/div/button
```

As you can see, the XPath selector leading to the button depends on the elements above it. So, if a web developer decides to change the layout and switches the div to something else, then our XPath—and our scraper—is broken.

I wanted a tool that would let me describe a scrape using a set of standardized options. After partnering with a few journalism organizations, I came up with three common investigative scrape patterns:

1. Crawl a site and download documents (e.g., download all PDFs on a site).

2. Query and grab a single result page (e.g., enter a package tracking number and get a page saying where it is).

3. Query and save all result pages (e.g., a Google search with many results pages).

The result of this work is AutoScrape, a resilient web scraper that can perform these kinds of scrapes and survive a variety of site changes. Using a set of basic configuration options, most common journalistic scraping problems can be solved. It operates on three main principles in order to reduce the amount of maintenance required when running scrapes:

**Navigate like people do.**

When people use websites they look for visual cues. These are generally standardized and rarely change. For example, search forms typically have a "Search" or "Submit" button. AutoScrape takes advantage of this. Instead of collecting DOM element selectors, users only need to provide the text of pages, buttons and forms to interact with while crawling a site.

**Extraction as a separate step.**

Web scrapers spend most of their time interacting with sites, submitting forms and following links. Piling data extraction onto that process increases the risk of having the scraper break. For this reason, AutoScrape separates site navigation and data extraction into two separate tasks. While AutoScrape is crawling and querying a site, it saves the rendered HTML for each page visited.

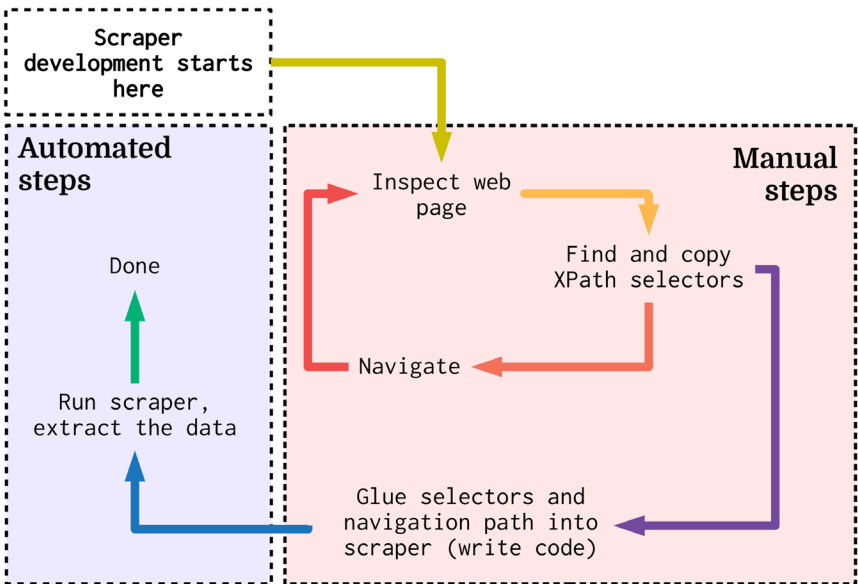**Protect extraction from page changes as much as possible.**

Extracting data from HTML pages typically involves taking XPaths, extracting data into columns, converting these columns of data into record rows and exporting them. In AutoScrape, we avoid this entirely by using an existing domain-specific template language for extracting JSON from HTML called Hext. Hext templates, as they're called, look a lot like HTML but include syntax for data extraction. To ease the construction of these templates, the AutoScrape system includes a Hext template builder. Users load one of the scraped pages containing data and, for a single record, click on and label each of the values in it. Once this is done, the annotated HTML of the selected record can be converted into a Hext template. JSON data extraction is then a matter of bulk processing the scraped pages with the Hext tool and template.

Hext templates are superior to the traditional XPath method of data extraction on pages that contain few class names or IDs, as is common on primitive government sites. While an XPath to an element can be broken by changes made to ancestor elements, breaking a Hext template requires changing the actual chunk of HTML that contains data.

# Using AutoScrape

Here are a few examples of what various scrape configurations look like using AutoScrape. In all of these cases, we're just going to use the command line version of the tool: scrape.py.

Let's say you want to crawl an entire website, saving all HTML and style sheets (no screenshots):

*A flowchart of the typical web scraper development process. Development starts at the top left and continues with manual steps: analyzing the website, reading the source code, extracting XPaths and pasting them into code. The single automated step consists of running the scraper. Because scrapers built this way are prone to breaking or extracting subtly incorrect data, the output needs to be checked after every scrape. This increases the overall maintenance requirements of operating a web scraper.*

```
./scrape.py \
  --maxdepth -1 \
  --output crawled_site \
  'https://some.page/to-crawl'
```

In the above case, we've set the scraper to crawl infinitely deep into the site. But if you want to only archive a single webpage, grabbing both the code and a full length screenshot (PNG) for future reference, you could do this:

```
./scrape.py \
  --full-page-screenshots \
  --load-images \
  --maxdepth 0 \
  --save-screenshots \
```

```
--driver Firefox \
--output archived_webpage \
'https://some.page/to-archive'
```

Finally, we have the real magic: interactively querying web search portals. In this example, we want AutoScrape to do a few things: load a webpage, look for a search form containing the text "SEARCH HERE", select a date (January 20, 1992) from a date picker, enter "Brandon" into an input box and then click "NEXT ->" buttons on the result pages. Such a scrape is described like this:

```
./scrape.py \
  --output search_query_data \
  --form-match "SEARCH HERE" \
  --input "i:0:Brandon,d:1:1992-01-20" \
  --next-match "Next ->" \
  'https://some.page/search?s=newquery'
```

All of these commands create a folder that contains the HTML pages encountered during the scrape. The pages are categorized by the general type of page they are: pages viewed when crawling, search form pages, search result pages and downloads.

```
autoscrape-data/
├── crawl_pages/
├── data_pages/
├── downloads/
└── search_pages/
```

Extracting data from these HTML pages is a matter of building a Hext template and then running the extraction tool. Hext templates can either be written from scratch or by using the Hext builder web tool included with AutoScrape. (https://github.com/brandonrobertz/autoscrape-py)

Currently, I'm working with the Computational Journalism Workbench team to integrate AutoScrape into their web platform so that you won't need to use the command line at all. Until that happens, you can go to the GitHub repo to learn how to set up a graphical version of AutoScrape.

# Fully Automated Scrapers and the Future

In addition to building a simple, straightforward tool for scraping websites, I had a secondary, more lofty motive for creating AutoScrape: using it as a testbed for fully automated scrapers.

Ultimately, I want to be able to hand AutoScrape a URL and have it automatically figure out what to do. Search forms contain clues about how to use them in both their source code and displayed text. This information can likely be used to by a machine learning[10] model to figure out how to search a form, opening up the possibility for fully-automated scrapers.

If you're interested in improving web scraping, or just want to chat, feel free to reach out. I'm in this for the long haul. So stay tuned.

# Glossary of Terminology

**1. Dataset** A collection of machine readable records, typically from a single source. A dataset can be a single file (Excel or CSV), a database table, or a collection of documents. In machine learning, a dataset is commonly called a corpus. When the dataset is being used to train[18] a machine learning model[12], it can be called a training dataset (a.k.a. a training set). Datasets need to be transformed into a matrix[11] before they can be used by a machine learning model.

**2. Distance Function, Distance Metric** A method for quantifying how dissimilar, or far apart, two records are. Euclidean distance, the simplest distance metric used, is attributed to the Ancient Greek mathematician Euclid. This distance metric finds the length of a straight line between two points, as if using a ruler. Cosine distance is another popular metric which measures the angle between two points using trigonometry.

**3. Document Object Model, DOM** A representation of a HTML page using a hierarchical tree. This is the way that browsers "see" web pages.

**4. Feature** A column in a dataset representing a specific type of values. A feature is typically represented as a variable in a machine learning model. For example, in a campaign finance dataset, a feature might be "contribution amount" or "candidate name". The number of features in a dataset determines its dimensionality. In many machine learning algorithms, high dimensional data (data with lots of features) is notoriously difficult to work with.

**5. Hash** A short label or string of characters identifying a piece of data. Hashes are generated by a hash function. An example of this

comes from the most common use case: password hashes. Instead of storing passwords in a database for anyone to read (and steal), password hashes are stored. For example, the password "Thing99" might get turned into something like b3aca92c793ee0e9b1a9b0a5f-5fc044e05140df3 by a hash function and saved in a database. When logging in, the website will hash the provided password and check it against the one in the database. A strong cryptographic hash function can't feasibly be reversed and uniquely identifies a record. In other usages, such as in_ LSH_, a hash may identify a group of similar records. Hashes are a fixed length, unlike the input data used to create them.

**6. *k-NN, k-Nearest Neighbors*** An algorithm for finding the k number of most similar records to a given record, or query point. k-NN can use a variety of distance metrics[2] to measure dissimilarity, or distance_, between points. In k-NN, when _k is equal to 1, the algorithm will return the single most similar record. When k is greater than 1, the algorithm will return multiple records. A common practice is to take the similar records, average them and make educated guesses about the query point.

**7. *Locality Sensitive Hashing, LSH*** A method, similar in application to k-NN[6], for identifying similar records given a query record. LSH uses some statistical tricks like hashing[5] and projection[9] to do this as a performance optimization. Due to this, it can be used on large amounts of data. The penalty for this is that it's possible for false records to turn up in the results and, inversely, for actual similar records to be missed.

**8. *Preprocessing*** A step in data analysis that happens before any actual analysis occurs to transform the data into a specific format or to clean it. A common preprocessing task is lowercasing and stripping symbols. Vectorization[20] is a common preprocessing step found in machine learning and statistics.

**9. *Projection*** A mathematical method for taking an input vector[20] and transforming it into another dimension. Typically, this is done by taking high dimensional data (data with a large number of columns) and converting it to a lower dimension. A simple example of

this would be taking a 3D coordinate and turning it into a 2D point. This is one of the key concepts behind LSH[7].

**10. Machine Learning, ML** A field of statistics and computer science focused on building or using algorithms that can perform tasks, without being told specifically how to accomplish them, by learning from data. The type of data required by the machine learning algorithm, labeled or unlabeled, splits the field into two major groups: supervised[17] and unsupervised[19], respectively. Machine learning is a subfield of Artificial Intelligence.

**11. Matrix** Rectangularly arranged data made up of rows and columns. In the machine learning context, every cell in a matrix is a number. The numbers in a matrix may represent a letter, number or category.

```
        Example m-by-n matrix (2x3)

     n columns (3)
    .--------------------------------.
  m | 11.347271944951725 | 2203 | 2.0 | <- row vector
rows |-------------------+------+-----|
 (2) | 9.411296351528783  | 1867 | 1.0 |
    `--------------------------------'
        \
         This is element (2,1)
```

Each row, which represents a single record, is known as a vector[20]. The process of turning source data into a matrix is known as vectorization.

**12. Model, Statistical Model** A collection of assumptions that a machine learning algorithm has learned from a dataset. Fundamentally, a model consists of numbers, known as weights, that can be be plugged into a machine learning algorithm. We use models to get data out of machine learning algorithms.

**13. Outlier Detection** A method for identifying records that are out of place in the context of a dataset. These outlying data points can be thought of as strange, suspicious, fraudulent, rare, unique, etc. Outlier detection is a major subfield of machine learning with applications in fraud detection, quality assurance and alert systems.

**14. *Regression*** A statistical method for identifying relationships among the features[4] in a dataset.

**15. *Scraping, Web Scraping*** The process of loading a web page, extracting information and collecting it into a specific structure (a database, spreadsheet, etc). Typically web scraping is done automatically with a program, or tool, known as a web scraper.

**16. *String*** A piece of data, arranged sequentially, made up of letters, numbers or symbols. Technically speaking, computers represent everything as numbers, but they are converted to letters when needed. Numbers, words, sentences, paragraphs and even entire documents can be represented as strings.

**17. *Supervised Machine Learning, Supervised Learning*** A subfield of machine learning where algorithms learn to predict values or categories from human-labeled data. Examples of supervised machine learning problems: (1) predicting the temperature of a future day using a dataset of historical weather readings and (2) classifying emails by whether or not they are spam from a set of categorized emails. The goal of supervised machine learning is to learn from one dataset and then make accurate predictions on new data (this is known as generalization).

**18. *Training*** The process of feeding a statistical or machine learning algorithm data for the purpose of learning to predict, identifying structure, or extracting knowledge. As an example, consider a list of legitimate campaign contributions. Once an algorithm has been shown this data, it generates a model[12] representing how these contributions typically look. This model can be used to spot unusual contributions, since the model has learned what normal ones look like. There are many different methods for training models, but most of them are iterative, step-based procedures that slowly improve over time. A common analogy for how models are trained is hill climbing: knowing that a flat area (a good solution) is at the top of a hill, but only being able to see a short distance due to thick fog, the top can be found by following steep paths. Training is also known as model fitting.

**19. *Unsupervised Machine Learning, Unsupervised Learning*** A subfield

of machine learning where algorithms learn to identify the structure or find patterns within a dataset. Unsupervised algorithms don't require human labeling or organization, and therefore can be used on a wide variety of datasets and in many situations. Examples of unsupervised use cases: (1) discovering natural groups of records in a dataset, (2) finding similar documents in a dataset and (3) identifying the way that events normally occur and using this to detect unusual events (a.k.a. outlier detection and anomaly detection[13]).

*20. Vectorization, Vector* The process of turning a raw source dataset into a numerical matrix[11]. Each record becomes a row of the matrix, known as a vector.

*21. Weight* A number that is used to either increase or decrease the importance of a feature[4]. Weights are used in supervised machine learning to quantify how well one variable predicts another; in unsupervised learning, weights are used to emphasize features that segment a dataset into groups.

*22. XPath* A description of the location of an element on a web page. From the browser's perspective, a web page is represented as a hierarchical tree known as the Document Object Model (DOM)[3]. An XPath selector describes a route through this tree that leads to a specific part of the page.

```
    _____
   /    /|
  /  /| |
 /  /_|_|
/_____/__/
      /   /
     /  /
    /  /
   /__/
```

artificialinformer.com