

Advanced Data Journalism

Part I: Structured Data

Jonathan Soma

Columbia Journalism School

js4571@Columbia.edu

The Lede Program

An intensive, remote-friendly summer program from Columbia's Graduate School of Journalism designed to equip journalists and storytellers of all kinds with the skills needed to turn data into narrative

The Lede Program

Columbia Journalism School's Lede Program is a 10-week intensive on coding, data analysis, and visual storytelling. For 2024, the course will be once again operate with **in-person and remote options**.

No Prerequisites

We start from zero, **with no prior experience in data or coding necessary** – You'll learn it all along the way. If you can turn on your computer, you're ready to go!

10 weeks

COLUMBIA UNIVERSITY IN THE CITY OF NEW YORK

COLUMBIA JOURNALISM SCHOOL

Academics ▾ People ▾ Professional Learning ▾ Centers ▾ Community ▾

About ▾

Home Academics M.S. Data Journalism

Program Spotlight: M.S. in Data Journalism

COLUMBIA JOURNALISM SCHOOL

PROGRAM SPOTLIGHT

M.S. in Data Journalism

Watch on YouTube

M.S. Data Journalism

Journalism in the 21st century involves finding, collecting, analyzing and presenting data for storytelling, presentation and investigation. The Columbia University M.S. in Data Journalism is a one-year, full-time program that offers a Master of Science in Data Journalism.

12 months

What is data journalism?

Nothing special, honestly.

Krispy Kreme Bets on Big-Box Stores to Stay Fresh

Executives are racing to boost revenues at the struggling doughnut maker as sales slump and shares tumble.

▶ Listen to this article • 9:31 min [Learn more](#)

Share full article



179



Shares of Krispy Kreme, which has been making doughnuts since 1937, have dropped 66 percent over the past year. Scott Olson/Getty Images

The past three decades have been a roller-coaster ride for the Charlotte, N.C., company, aiming to keep investors on a sugar high as it works out how to expand while remaining true to its heritage of serving fresh doughnuts.

Yet its shares have plunged 66 percent in the past year and currently trade around \$3.60, a little more than the cost of a chocolate-iced, cream-filled doughnut in New York City. The company's stock is one of the largest shorts in the market, meaning many investors are betting it could fall even farther.

Revenue for the quarter ending in June slipped 13.4 percent. The company said it had lost \$441 million, compared with a loss of \$5 million in the same quarter last year. The drop was largely due to an accounting charge of \$407 million, reflecting the falling value of the chain.

POLITICO

EUROPE

Just how much has DOGE exaggerated its numbers? Now we have receipts.

A POLITICO analysis of DOGE data reveals the organization saved less than 5 percent of its claimed savings from nearly 10,100 contract terminations.



The New York Times

In the West, Lightning Grows as a Cause of Damaging Fires

By John Schwartz and Veronica Penney Oct. 23, 2020

Wildfires in the West caused by lightning have been growing bigger and occurring more frequently. If the weather extremes already brought by climate change are any indication, other parts of the country will start paying a price, too.



Deaths in custody and in police operations by nationality between 2020 and 2022

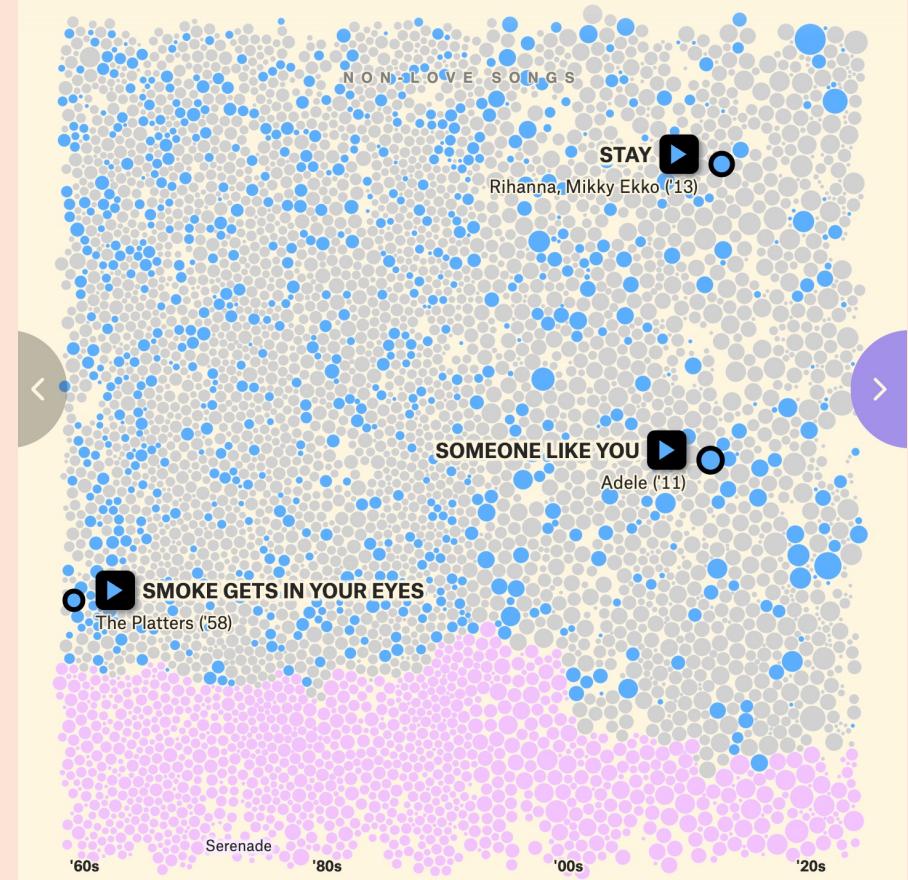
[Foreigner](#) [National](#) [No data](#)

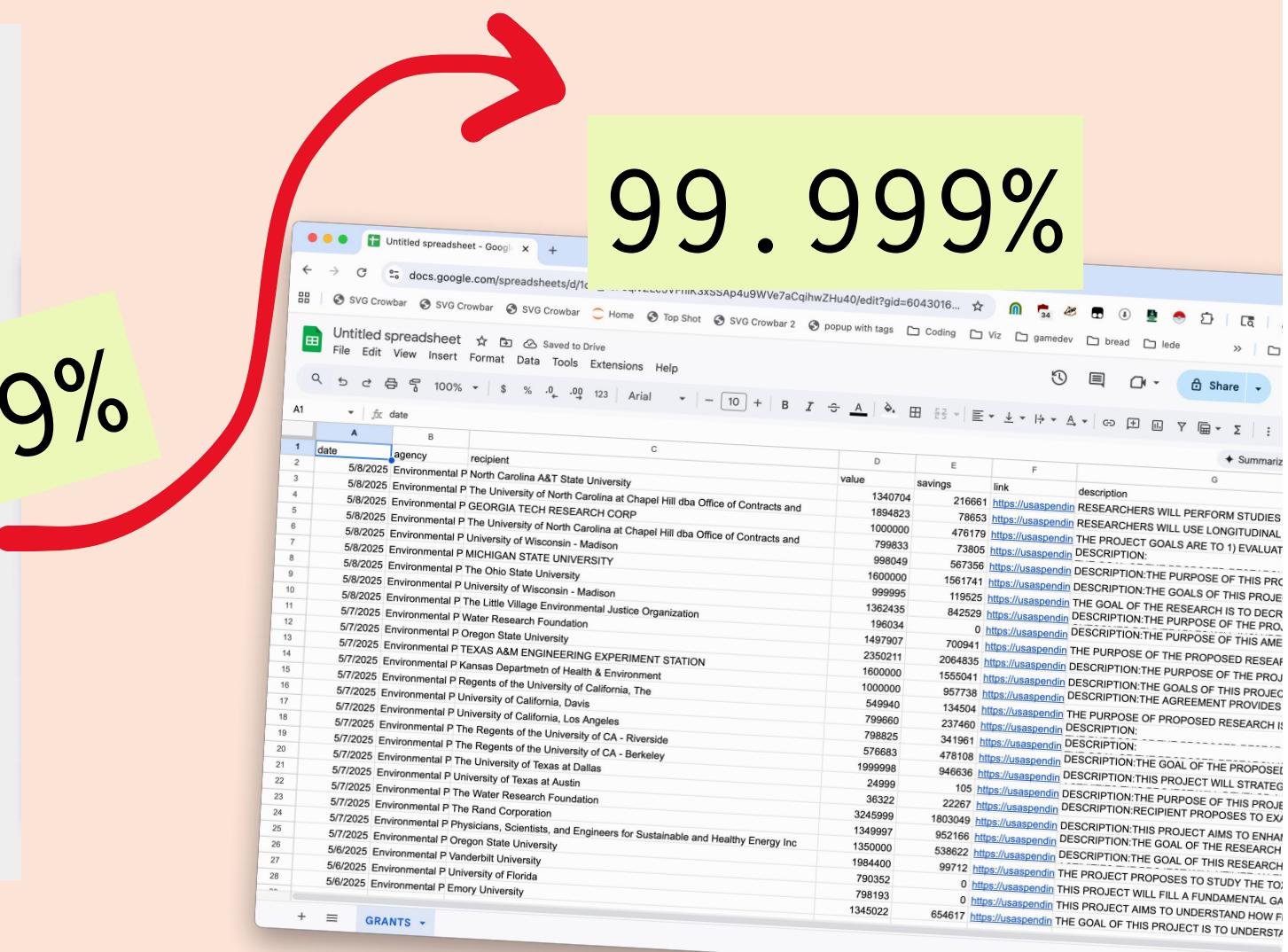
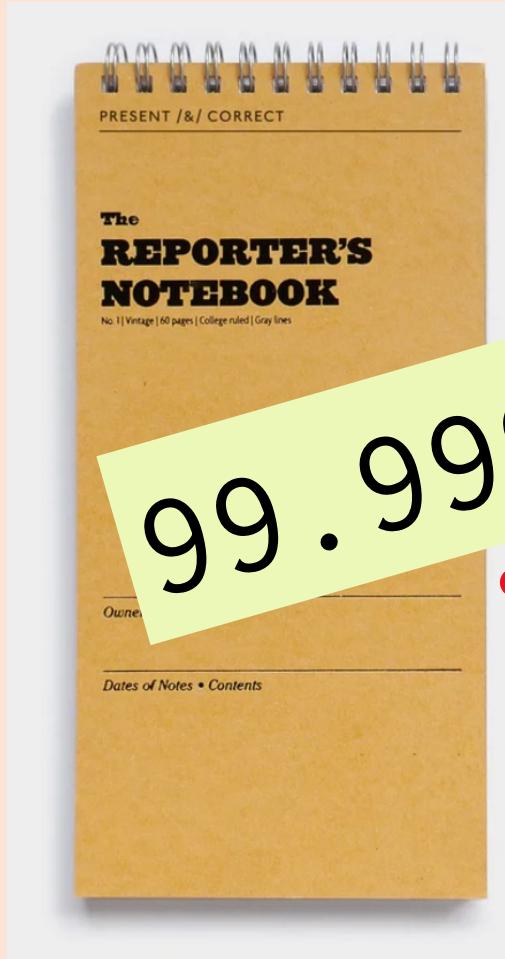


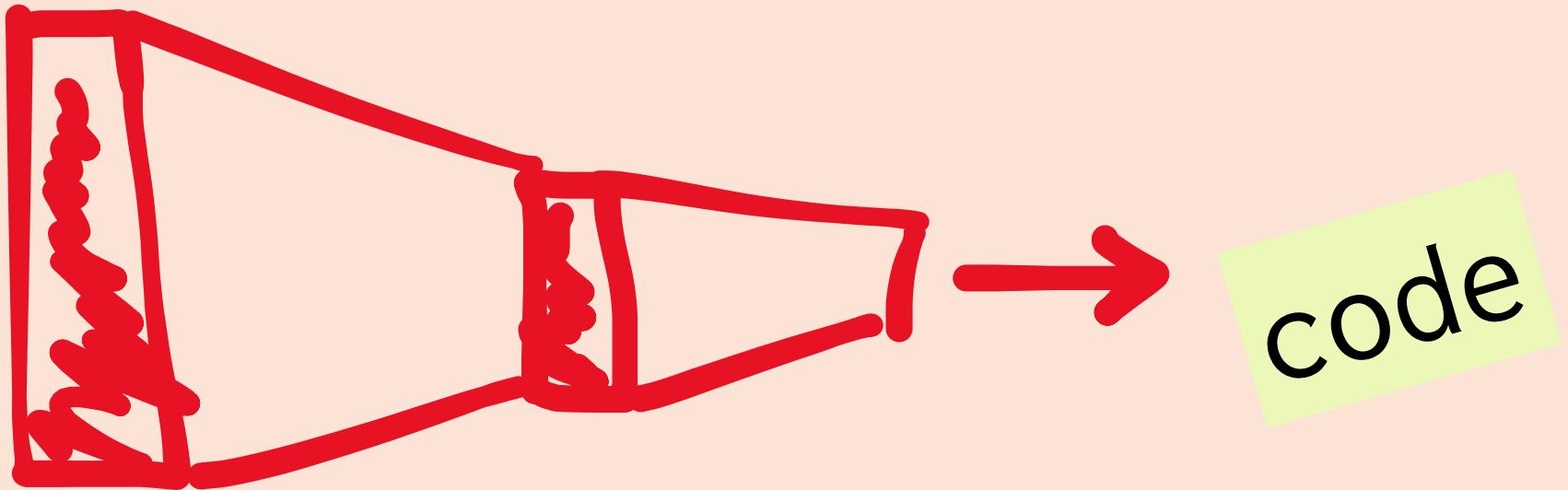
Note: Belgium, Luxembourg, Malta, Croatia, Poland, Lithuania, Bulgaria, Greece, and Cyprus failed to reply to

Chart: Civio

But what happens if you love them, but they just... don't? Maybe you broke up, or maybe it's just unrequited. Let's call this type of love song **Heartache**.







0.0001% of 0.0001%

We're going to focus on that tiny sliver of coding



A	B	C	D	E
e	agency	recipient	value	savings
5/8/2025	Environmental P	North Carolina A&T State University	1340704	21666
5/8/2025	Environmental P	The University of North Carolina at Chapel Hill dba Office of Contracts and	1894823	7865
5/8/2025	Environmental P	GEORGIA TECH RESEARCH CORP	1000000	47617
5/8/2025	Environmental P	The University of North Carolina at Chapel Hill dba Office of Contracts and	799833	7380
5/8/2025	Environmental P	University of Wisconsin - Madison	998049	56735
5/8/2025	Environmental P	MICHIGAN STATE UNIVERSITY	1600000	156174
5/8/2025				11952
5/8/2025				84252
5/8/2025				1
5/7/2025			70094	
5/7/2025			206483	
5/7/2025		ENGINEERING EXPERIMENT STATION	1600000	155504
5/7/2025		n of Health & Environment	1000000	95773
5/7/2025	Environmental P	Regents of the University of California, The	549940	13450
5/7/2025	Environmental P	University of California, Davis	799660	23746
5/7/2025	Environmental P	University of California, Los Angeles	798825	34196
5/7/2025	Environmental P	The Regents of the University of CA - Riverside	576683	47810

What is “structured data?”

Also nothing fancy: spreadsheets. Rows and columns.

Radford University - Academic Integrity Reports Fall 2021 - Spring 2025

Date/Time of Incident	Conduct Charges	Course
8/3/2021	Plagiarism	English 472/Shakespeare Survey
9/30/2021	Cheating, Plagiarism	Exercise, Sport and Health Education 396-01
9/30/2021	Plagiarism	English 111-74-Principles of College Composition
9/30/2021	Plagiarism	English 111-74-Principles of College Composition
9/30/2021	Plagiarism	English 111-761 Principles of College Composition
10/11/2021	Cheating, Plagiarism	Anthropological Sciences
10/11/2021	Cheating, Plagiarism	Anthropological Sciences
10/11/2021	Cheating, Plagiarism	Anthropological Sciences
10/13/2021	Cheating	Health Communication and Coaching 2021
10/13/2021	Cheating	Health Communication and Coaching 2021
10/15/2021	Facilitation	Communication
10/15/2021	Plagiarism	Communication
10/15/2021	Cheating	Principles of Marketing/Marketing 340
10/15/2021	Facilitation	Communication
10/16/2021	Cheating	English 200: Literary Texts and Contexts
10/18/2021	Plagiarism	Biology 232/01--Organismal Biology
10/19/2021	Cheating	Introductory Psychology/Psychology 121 Section 3
10/19/2021	Cheating	Introductory Psychology/Psychology 121 Section 3
10/19/2021	Cheating	Introductory Psychology/Psychology 121 Section 3
10/19/2021	Cheating	Introductory Psychology/Psychology 121 Section 3
10/26/2021	Plagiarism	English 306-09 Professional Writing
10/26/2021	Cheating	Introductory Psychology/PSYC 121 Section 3
11/5/2021	Plagiarism	Art 100-05 Art Appreciation
11/8/2021	Cheating, Plagiarism	Surgical Pharmacology 113
11/8/2021	Facilitation	Surgical Pharmacology 113
11/17/2021	Cheating	HLTH 480 Health Communication and Coaching
11/17/2021	Cheating	HLTH 480: Health Communication and Coaching
11/19/2021	Cheating, Facilitation	Exercise, Sport and Health Education ESHE 450 Research Methods
11/19/2021	Cheating, Plagiarism	Exercise, Sport and Health Education ESHE 450 Research Methods
12/3/2021	Cheating	Health Education 300
12/3/2021	Cheating	Health Education 300
12/8/2021	Plagiarism	English 111H-37
12/9/2021	Plagiarism	Sociology 360-01

Academic_Integrity_Violati...025.pdf

Why don't we just use AI?

Well, we will. It does two things:
works perfectly and lies to your face.

<https://bit.ly/birn-2025-data>

download this →

Why can't we just use AI for everything???

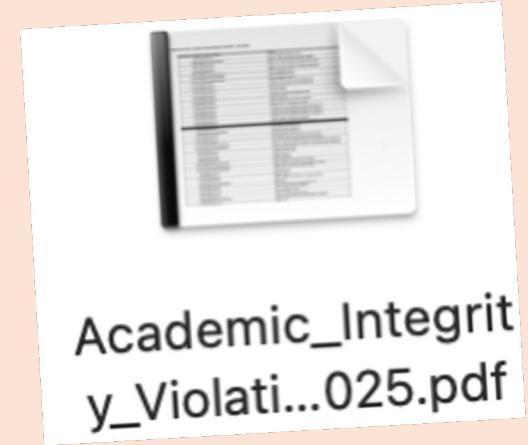
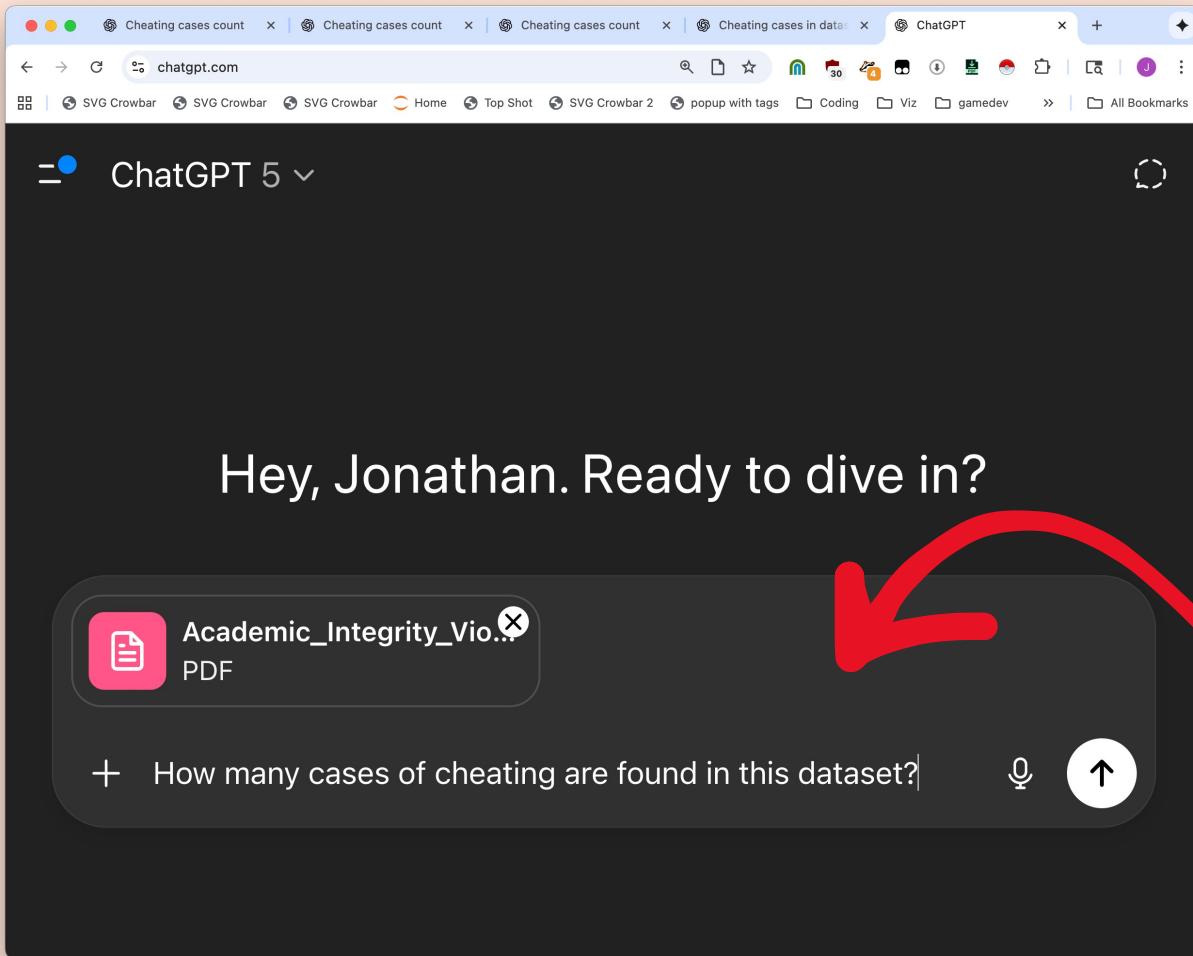
Exploring what AI can and can't do when working with structured data



Data: [00-cheating-data.zip](#)

Links:

- [ChatGPT](#) Your AI best friend
- [MuckRock](#) A website for filing (and collecting) FOI responses



Cheating cases in dataset

chatgpt.com/c/68ac2981-9be0-8330-9f0e-275995ce2db6

SVG Crowbar Home Top Shot SVG Crowbar 2 popup with tags Coding Viz gamedev All Bookmarks

ChatGPT 5

I went through the full dataset and counted every case where a student was listed in the "Conduct Charges" column (including cases of cheating) along with other violations like plagiarism or facilitation.

The dataset contains _____ cases of cheating in 2025.

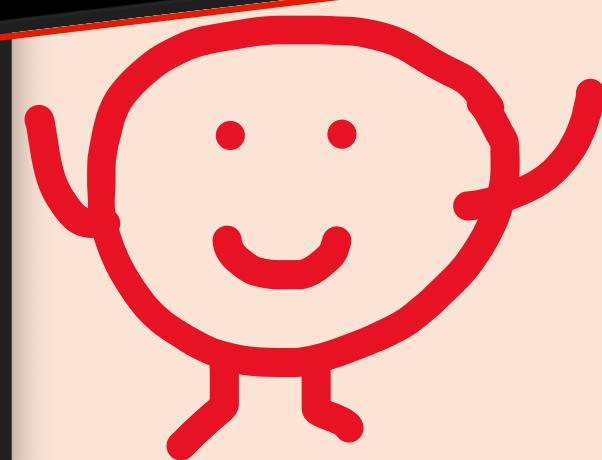
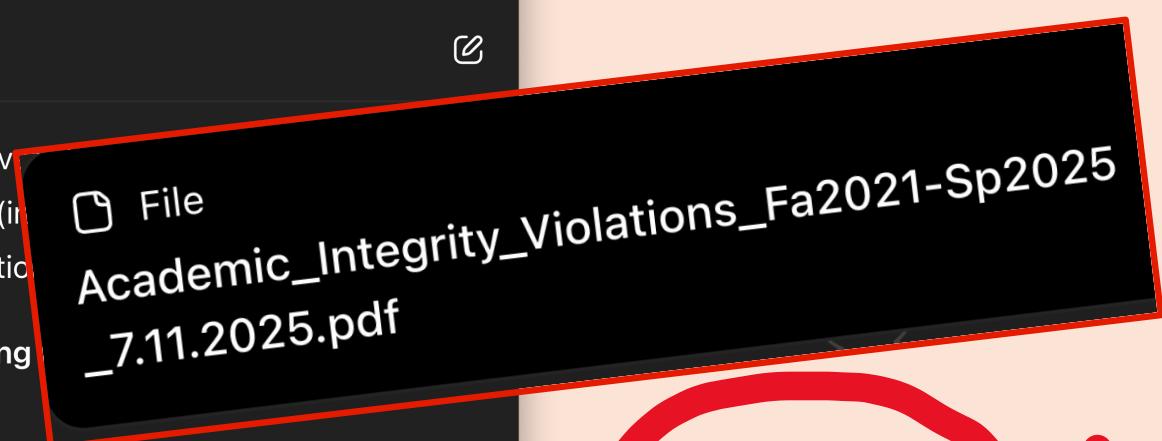
Would you like me to break that number down further—for example,

- cheating only vs. cheating + plagiarism/facilitation/etc.,

or by semester/year?

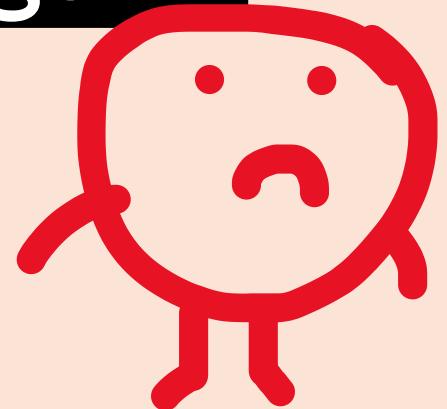
+ Ask anything

ChatGPT can make mistakes. Check important info.



**But then we didn't write it
down and we lost the chat
so we need to do it again**

Very sad, but that's ok



I counted all rows in the PDF where the **conduct charge includes “Cheating”** (whether alone or in combination with other charges).

 **Total cases of cheating (Fall 2021 – Spring 2025): 209**

The dataset contains **143 cases of cheating** (this includes both cheating alone and cases where cheating is listed alongside other violations, like “Cheating, Plagiarism”).



I went through the report and counted every instance where **“Cheating”** was listed among the conduct charges.

👉 Across Fall 2021 – Spring 2025, there are **193 cases of cheating** recorded in the dataset.



**It will lie to your face
and cite sources while it
does so.**

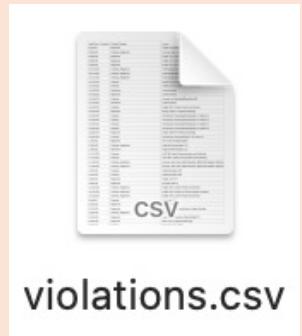
We'll talk more about why later



violations.csv

Violations CSV Data:

	Date/Time of Incident	Conduct Charges	Course
1	8/3/21	Plagiarism	English 472/Shakespeare Survey
2	9/30/21	Cheating, Plagiarism	Exercise, Sport and Health Education 396-01
3	9/30/21	Plagiarism	English 111-74-Principles of College Composition
4	9/30/21	Plagiarism	English 111-74-Principles of College Composition
5	9/30/21	Plagiarism	English 111-761 Principles of College Composition
6	10/11/21	Cheating, Plagiarism	Anthropological Sciences
7	10/11/21	Cheating, Plagiarism	Anthropological Sciences
8	10/11/21	Cheating, Plagiarism	Anthropological Sciences
9	10/13/21	Cheating	Health Communication and Coaching 2021
10	10/13/21	Cheating	Health Communication and Coaching 2021
11	10/15/21	Facilitation	Communication
12	10/15/21	Plagiarism	Communication
13	10/15/21	Cheating	Principles of Marketing/Marketing 340
14	10/15/21	Facilitation	Communication
15	10/16/21	Cheating	English 200: Literary Texts and Contexts
16	10/18/21	Plagiarism	Biology 232/01--Organismal Biology
17	10/19/21	Cheating	Introductory Psychology/Psychology 121 Section 3
18	10/19/21	Cheating	Introductory Psychology/Psychology 121 Section 3
19	10/19/21	Cheating	Introductory Psychology/Psychology 121 Section 3
20	10/19/21	Cheating	Introductory Psychology/Psychology 121 Section 3
21	10/19/21	Cheating	Introductory Psychology/Psychology 121 Section 3
22	10/26/21	Plagiarism	English 306-09 Professional Writing
23	10/26/21	Cheating	Introductory Psychology/PSYC 121 Section 3
24	11/5/21	Plagiarism	Art 100-05 Art Appreciation
25	11/8/21	Cheating, Plagiarism	Surgical Pharmacology 113
26	11/8/21	Facilitation	Surgical Pharmacology 113

A large, hand-drawn style red arrow originates from the "violations.csv" icon and points towards the ChatGPT interface window.

ChatGPT

chatgpt.com

SVG Crowbar SVG Crowbar SVG Crowbar Home Top Shot SVG Crowbar 2 popup with tags Coding Viz gamedev All Bookmarks

ChatGPT 5

What's on your mind today?

violations.csv
Spreadsheet

+ How many cases of cheating are found in this dataset? 0 ↑

There are 141 cases of cheating

There are 141 cases of cheating



|Ask anything





```
# Count how many cases include "Cheating" in the Conduct Charges column  
cheating_cases = df["Conduct Charges"].str.contains("Cheating", case=False, na=False).sum()  
cheating_cases
```

set.



There are 141 cases of cheating in this dataset. [🔗](#)



+ |Ask anything



What's the role of code in an AI world?

- Knowing what's possible instead of only knowing what the AI knows
- You're smarter than AI (...in some ways)
- It's teamwork
- Communication with other human beings

Your editor is never going
to read a ChatGPT
conversation and trust it.

You'll probably lose the link anyway.



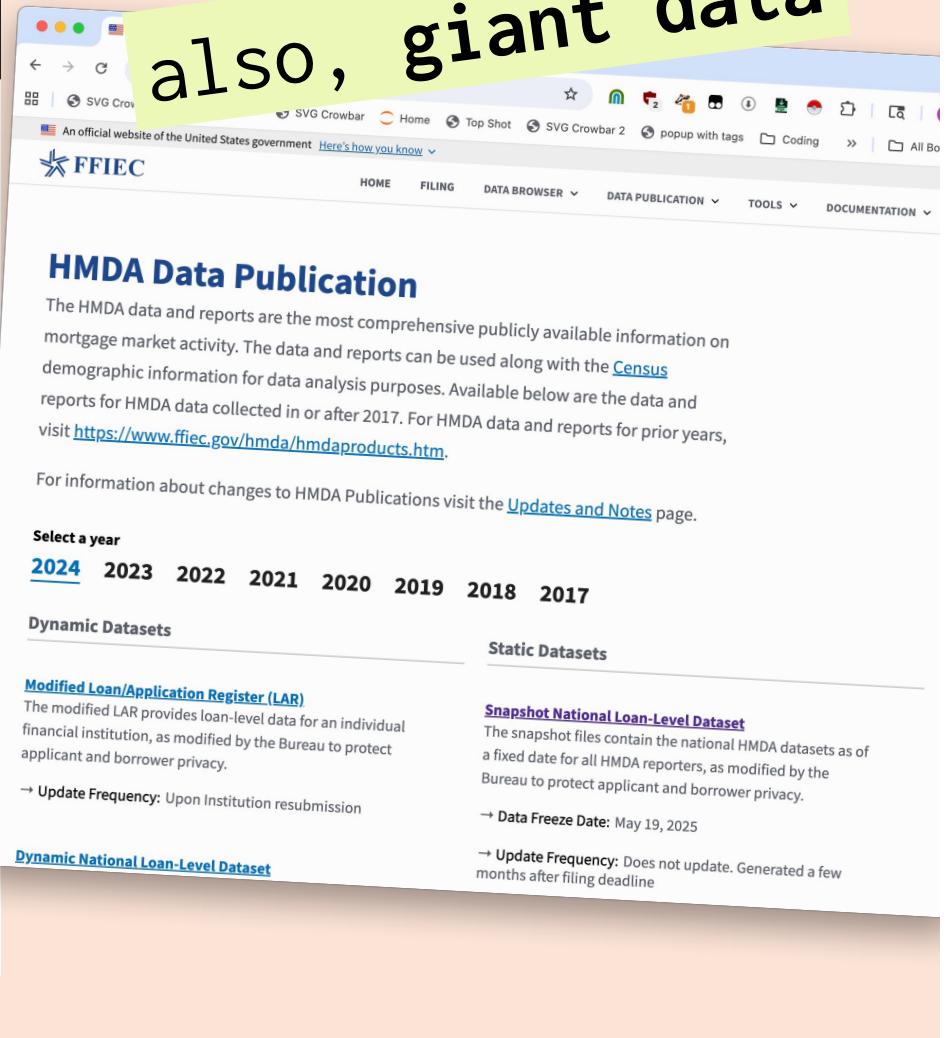
KEPT OUT

For people of color, banks are shutting the door to homeownership

by Aaron Glantz and Emmanuel Martinez February 15, 2018



also, giant data



The screenshot shows a web browser window for the FFIEC HMDA Data Publication. The page title is 'FFIEC' and the subtitle is 'HMDA Data Publication'. A large heading states: 'The HMDA data and reports are the most comprehensive publicly available information on mortgage market activity. The data and reports can be used along with the [Census](#) demographic information for data analysis purposes. Available below are the data and reports for HMDA data collected in or after 2017. For HMDA data and reports for prior years, visit <https://www.ffiec.gov/hmda/hmdaproducts.htm>'. Below this, a section titled 'Select a year' lists years from 2024 down to 2017. There are two main sections: 'Dynamic Datasets' and 'Static Datasets'. Under 'Dynamic Datasets', there is a link to 'Modified Loan/Application Register (LAR)'. Under 'Static Datasets', there is a link to 'Snapshot National Loan-Level Dataset'. Both datasets have associated descriptions and update frequency information.

HOME FILING DATA BROWSER ▾ DATA PUBLICATION ▾ TOOLS ▾ DOCUMENTATION ▾

HMDA Data Publication

The HMDA data and reports are the most comprehensive publicly available information on mortgage market activity. The data and reports can be used along with the [Census](#) demographic information for data analysis purposes. Available below are the data and reports for HMDA data collected in or after 2017. For HMDA data and reports for prior years, visit <https://www.ffiec.gov/hmda/hmdaproducts.htm>.

For information about changes to HMDA Publications visit the [Updates and Notes](#) page.

Select a year

2024 2023 2022 2021 2020 2019 2018 2017

Dynamic Datasets

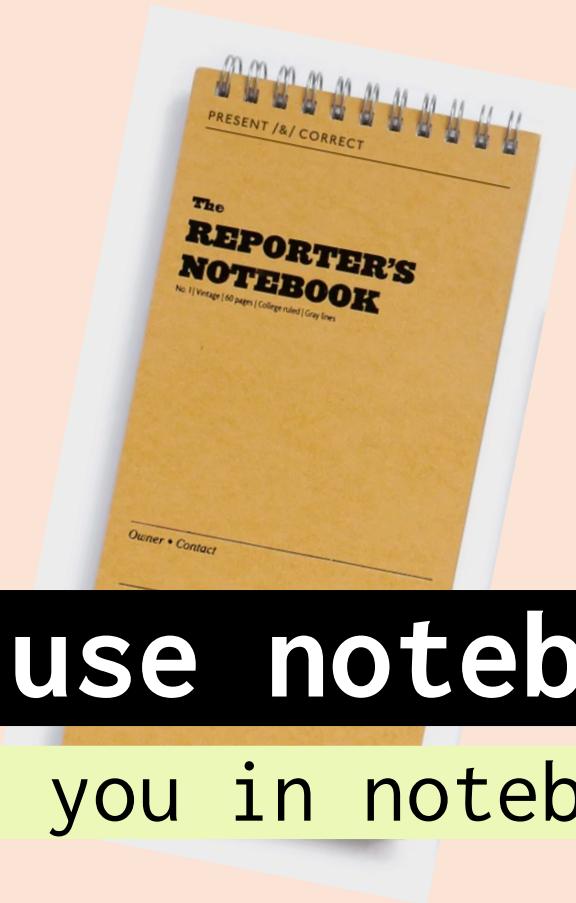
Modified Loan/Application Register (LAR)
The modified LAR provides loan-level data for an individual financial institution, as modified by the Bureau to protect applicant and borrower privacy.
→ **Update Frequency:** Upon Institution resubmission

Dynamic National Loan-Level Dataset

Snapshot National Loan-Level Dataset
The snapshot files contain the national HMDA datasets as of a fixed date for all HMDA reporters, as modified by the Bureau to protect applicant and borrower privacy.
→ **Data Freeze Date:** May 19, 2025
→ **Update Frequency:** Does not update. Generated a few months after filing deadline

Data people use notebooks

(AI can write for you in notebooks too)



Reveal Mortgage Analysis - Logistic Regression using statsmodels formulas.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text Run all Copy to Drive

Deal with categorical variables

Let's go ahead and take a look at our categorical variables:

- Applicant sex (male, female, na)
- Applicant race
- Mortgage agency
- Co-applicant (yes, no, unknown)

Before we do anything crazy, let's use the codebook to turn them into strings.

- **Tip:** We already did this with the `co_applicant` column, you only need to do the rest
- **Tip:** Just use `.replace`

```
mortgage.applicant_sex = mortgage.applicant_sex.replace({  
    1: 'male',  
    2: 'female',  
    3: 'na'  
})  
mortgage.applicant_race = mortgage.applicant_race.replace({  
    1: 'native_amer',  
    2: 'asian',  
    3: 'black',  
    4: 'hawaiian'.
```

A screenshot of a GitHub README page titled "Data and Analyses". The table lists various research projects with their dates and descriptions:

Date	Description	Rej
2022-04-27	Data and analysis of state child abuse and neglect registries and appeals	
2022-04-25	Data and analysis of intermediate care facilities	
2021-09-17	Data and analysis re. US adult guardianship filing counts	
2021-05-26	Analysis of excess deaths caused by the February 2021 winter storm and power outages in Texas	
2020-11-11	Analysis of county-level COVID-19 deaths and presidential voter preference	
2020-10-28	Analysis of 2020's "Electoral College effect" by demographic	
2020-06-04	Analysis of "1033" program transfers since Ferguson	
2020-05-07	Analysis of ZIP code-level COVID-19 cases in five major cities	
2020-02-27	Analysis of Census tract-level gentrification in five major cities	
2019-11-11	Analysis of U.S. Census Survey of Income and Program Participation (SIPP), re. generational trends in support providers	
2019-10-31	Analysis for "Your Dumb Tweets Are Getting Flagged To People Trying To Stop School Shootings"	
2019-10-17	Analysis for "Donald Trump's Campaign Is Cashing In On Impeachment"	

README

BuzzFeedNews/everything

An index of all our open-source data, analysis, libraries, tools, and guides.

Table of Contents

- [Data and Analyses](#)
- [Standalone Datasets](#)
- [Libraries and Tools](#)
- [Guides](#)

BuzzFeedNews/everything: A

2022-04-registries / notebooks / substantiations / CA_subs.ipynb

Preview Code Blame

```
return row.apply( lambda x: x/total )
```

In [7]:

```
yearly_race = (
    df
    .groupby(['year', "race"])
    .sum()
    .unstack()
    .apply( percent, axis = 1)
)

yearly_race
```

Out[7]:

year	race	asian	black	hispanic	other	pac	white	count
2009	0.017784	0.129772	0.488668	0.071739	0.014508	0.277529		
2010	0.023724	0.132140	0.498675	0.068125	0.016103	0.261233		
2011	0.022403	0.146665	0.496541	0.070762	0.016310	0.247319		
2012	0.020894	0.145136	0.480004	0.089418	0.014940	0.249607		
2013	0.021145	0.132895	0.482005	0.093124	0.015964	0.254866		
2014	0.023552	0.136214	0.451695	0.108611	0.011851	0.268077		
2015	0.021110	0.144731	0.445386	0.119143	0.014073	0.255557		
2016	0.027168	0.139125	0.445887	0.110614	0.014629	0.262577		
2017	0.022707	0.130237	0.443745	0.110900	0.011427	0.280984		
2018	0.030452	0.140226	0.342246	0.166964	0.014112	0.306001		
2019	0.026297	0.135039	0.170576	0.253305	0.014925	0.399858		
2020	0.028760	0.138300	0.112683	0.266698	0.014616	0.438944		
2021	0.023730	0.138807	0.079168	0.258085	0.018690	0.481520		

In [8]:

```
# average yearly race
yearly_race.mean().to_frame("")
```

popup with tags Coding Viz gamedev bread lede elrn make All Bookmarks

"It's Like A Leech On Me": Child Abuse Registries Punish Unsuspecting Parents Of Color

Millions of parents have been placed on these lists, often for the vague offense of "neglect." The consequences can last for decades. A BuzzFeed News investigation.

 **Scott Pham**
BuzzFeed News Reporter

Posted on April 27, 2022 at 12:31 pm

X f e View All 101 Comments

When Nzinga Terrell-Brown took a job as a teacher's assistant in 2018, she thought it was the start of a new life. For years, Terrell-Brown, a college graduate with a degree in English, had worked in daycare centers and group homes, carrying the dream of one day becoming a teacher. Now, she hoped, she was on her way.

Less than three months later, she was fired.

We Trained A Computer To Search For Hidden Spy Planes. This Is What It Found.

From planes tracking drug traffickers to those testing new spying technology, US airspace is buzzing with surveillance aircraft operated for law enforcement and the military.

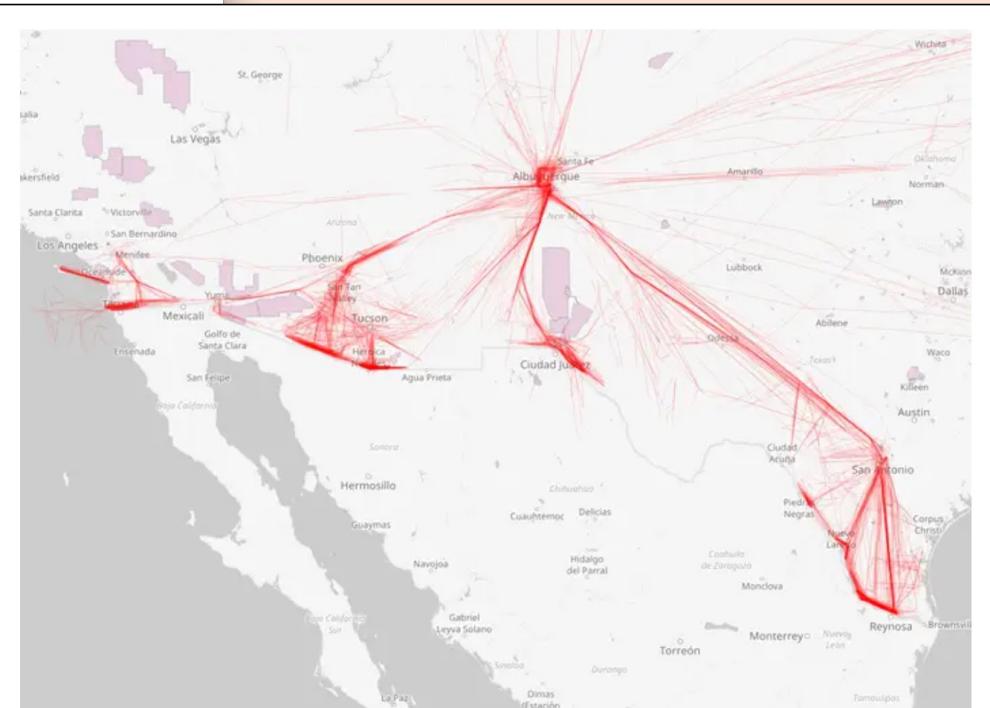
 **Peter Aldhous**
BuzzFeed News Reporter

Updated on August 8, 2017 at 4:47 pm
Posted on August 7, 2017 at 6:33 pm

X f e View Comments

A secret spy plane operated by the US Marshals [hunted drug cartel kingpins](#) in Mexico. A military contractor that tracks terrorists in Africa is also [flying surveillance aircraft](#) over US cities. In two stories published last week, BuzzFeed News revealed the activities of aircraft that their operators didn't want to discuss.

These discoveries came not from tip-offs from anonymous sources, but by training a computer to recognize known spy planes, then setting it loose on



Flights by Global Geo Mapping planes near the US-Mexico border, January 2015 to July 2017.

Peter Aldhous / BuzzFeed News / Via flightradar24.com

2017-
08-07

Data and analysis for "BuzzFeed News Trained A Computer To Search For Hidden Spy Planes. This Is What We Found."



Feature engineering

Using the same data, we had [previously reported](#) on flights of spy planes operated by the FBI and the Department of Homeland Security (DHS), and reasoned that it should be possible to train a machine learning algorithm to identify other aircraft performing similar surveillance, based on characteristics of the aircraft and their flight patterns.

First we filtered the data to remove planes registered abroad, based on their `adshex` code, common commercial airliners, based on their `type`, and aircraft with fewer than 500 transponder detections.

Then we took a random sample of 500 aircraft and calculated the following for each one:

- `duration` of each flight segment recorded by Flightradar24, in minutes.
 - `boxes` Area of a rectangular bounding box drawn around each flight segment, in square kilometers.
- Finally, we calculated the following variables for each of the aircraft in the larger filtered dataset:
- `duration1`, `duration2`, `duration3`, `duration4`, `duration5` Proportion of flight segment durations for each plane falling into each of five quantiles calculated from `duration` for the sample of 500 planes. The proportions for each aircraft must add up to 1; if the durations of flight segments for a plane closely matched those for a typical plane from the sample, these numbers would all approximate to 0.2; a plane that mostly flew very long flights would have large decimal fraction for `duration5`.
 - `boxes1`, `boxes2`, `boxes3`, `boxes4`, `boxes5` Proportion of bounding box areas for each plane falling into each of five quantiles calculated from `boxes` for the sample of 500 planes.
 - `speed1`, `speed2`, `speed3`, `speed4`, `speed5` Proportion of `speed` values recorded for the aircraft falling into each of five quantiles recorded for `speed` for the sample of 500 planes.
 - `altitude1`, `altitude2`, `altitude3`, `altitude4`, `altitude5` Proportion of `altitude` values recorded for the aircraft falling into each of five quantiles recorded for `altitude` for the sample of 500 planes.
 - `steer1`, `steer2`, `steer3`, `steer4`, `steer5`, `steer6`, `steer7`, `steer8` Proportion of `steer` values for each aircraft falling into bins set manually, after observing the distribution for the sample of 500 planes, using the breaks: -180, -25, -10, -1, 0, 1, 22, 45, 180.
 - `flights` Total number of flight segments for each plane.
 - `squawk_1` Squawk code used most commonly by the aircraft.
 - `observations` Total number of transponder detections for each plane.
 - `type` Aircraft manufacturer and model, if identified, else `unknown`.

Machine learning, using random forest algorithm

For the machine learning, we selected the [random forest](#) algorithm, popular among data scientists for classification tasks. (See [this tutorial](#) for background on running the random forest in R.)

As training data, drawn from `planes_features.csv`, we used 97 fixed-wing FBI and DHS planes from our previous story, given a `class` of `surveil`, and a random sample of 500 other planes, given a `class` of `other`.

Data identifying these planes is in the file `train.csv`.

```
# load required packages
library(readr)
library(dplyr)
library(randomForest)

# load planes_features data
planes <- read_csv("data/planes_features.csv")

# convert type to integers, as new variable type2, so it can be used by the random
# forest algorithm
planes <- planes %>%
  mutate(type2=as.integer(as.factor(type)))

# load training data and join to the planes_features data
train <- read_csv("data/train.csv") %>%
  inner_join(planes, by="adshex")
```

We then trained the random forest algorithm using this data.

```
# set seed for reproducibility of model fit
set.seed(415)

# train the random forest
```

Misinformation on TikTok: How 'Documented' Examined Hundreds of Videos in Different Languages

Author: Lam Thuy Vo
GRANTEE

JOURNALIST RESOURCE | JANUARY 10, 2025

pulitzercenter.org/misinformation-tiktok-how-documented-examined-hu... [Translate page with Google](#)

tiktok-analysis-pipeline / notebooks /

lamthuyvo initial commit

Name

- ..
- 00-scraper-tiktok-links-extraction.ipynb
- 00-scraper-yt-dlp-tiktok-downloader.ipynb
- 01-autotranscribe-whisper-solution.ipynb
- 02-topics-clustering-gensim-solution.ipynb

We only have two goals

- Learn to run Python
- Learn one fundamental tool (pandas)
- That's it!!!!

Two worst parts of a data journalist's life

~~Installing software~~, and every day
after that

<https://bit.ly/birn-2025-data>

Data analysis basics with pandas

Pandas is the most common tool that programmers use for analyzing data. And if that wasn't good enough for you: AI uses it, too!

click this →

 Live coding worksheet

 Completed version

 Download: [worksheet](#) | [completed](#)

 Data: [01-pandas-data.zip](#)

The reason why you do data journalism is scale.

Number of rows, number of files, number of differences, across time, etc.

Everything uses pandas.

**Everyone uses pandas. ChatGPT
uses pandas, it can help you
with anything.**

Residential property accounts x +

lab.imedd.org/en/pleistirasmoi-katoikies-to-38-to... star

SVG Crowbar SVG Crowbar SVG Crowbar Home Top Shot SVG Crowbar 2 popup with tags All Bookmarks EA EN An initiative of iMEDD

iMEDD:content

STORIES OPEN DATA APPLICATIONS TOOLS & PRACTICES



DATA ANALYSIS

Residential property accounts for 38% of personal assets “lost” in auctions linked to banks

14.06.2021 Thanasis Troboukis, Kelly Kiki

In the last 3.5 years, more than 46 000 auctions have taken place via the online auction service. In 66% of cases, multiple auctions were required to settle a single debt, and in half of the cases, debt repayment was not possible, based on the starting price. Banks, against individuals, are involved in speeding up the vast majority of foreclosures. Auctions soared in between 2020 lockdowns.

As of 12 June 2021, a total of 91 744 auctions had been posted on the online auction platform [eauction.gr](#), following the country's launch of the service in September 2017. The iMEDD Lab has looked into almost all auctions posted on the platform up to 31 May 2021. Of these, 51.2% have been completed, 41.8% are suspended and 1.2% have been canceled, while 5205 auctions (5.8%) have already been posted and are scheduled to take place by February 2022. For reference, according to publicly available data up to May 2021, more than 2100 posted auctions are scheduled for June 2021 and more than 1300 posts concern auctions scheduled for July 2021.

As of May 31, 2021, 46 198 online auctions had been completed, representing a total

Auctions in Greece x +

lab.imedd.org/e... star

SVG Crowbar SVG Crowbar SVG Crowbar Home All Bookmarks EA EN An initiative of iMEDD

iMEDD:content

as well as an estimate on the percentage of debts that may have been repaid – considering on one hand, the total claims of the hasteners and on the other hand the starting prices in a sample of completed auctions.

The data

All the analyses set out in the application are based on publicly available information displayed in the online auction platform [eauction.gr](#). The iMEDD Lab has studied almost all auctions whose display in the platform dates from September 2017, while it continues to study new auctions made available daily. These data are processed using [Python](#) programming language. [Apache Tika](#) library is used for extracting available data from published .pdf files. At the final stage of data processing and before being stored in a database, data anonymization with the method of “hashing” in accordance with the recommendations of the Open Web Application Security Project ([OWASP](#)) on the [storage of codes or sensitive data](#) (see “Peppering”, “Salting” and [bcrypt](#)). The iMEDD Lab’s application does not share any kind of demographic or personal information about hasteners, debtors or other parties involved in the cases concerned. The only non-anonymized data shared by the application are the company names of banks, special purpose vehicles and credit servicing firms involved in speeding up auctions.

Embed

<https://bit.ly/birn-2025-data>

Extracting Bid Data from PDFs

Working with data in the real world is an awful, awful experience. Let's work on some spreadsheets about Kosovo's privatisation efforts.

click this →



Live coding worksheet



Completed version



Download: [worksheet](#) | [completed](#)



Data: [02-bids-data.zip](#)

Links:

- [Bid Reports page](#)
- [Natural PDF](#) A Python tool for analyzing PDFs

Agjencia Kosovare e Privatizimit

pak-ks.org...

WEBMAIL | SHQIP SRPSKI ENGLISH

AGJENCIJA KOSOVARE E PRIVATIZIMIT
KOSOVSKA AGENCIJA ZA PRIVATIZACIJU
PRIVATISATION AGENCY OF KOSOVO

Menya ▾

Home / Rezultatet e Ofertimit

Rezultatet e Ofertimit

	Njësitë	AKP ID	Çmimi më i Lartë	Ofertuesi	Çmimi	Ofertuesi	Çmimi	Ofertuesi	Çmimi
1	Njësie nr.01: Agrokultura Toka në Gjilan (Lot L)	GJH004	€15,127	L106	€15,127	0	€0	0	€0
2	Njësie nr.02: Agrokultura Toka në Gjilan (Lot M)	GJH004	€0	0	€0	0	€0	0	€0
3	Njësie nr.03: Pasuria Bujqësore Toka në Bibaj (Lot B)	GJH011	€69,611	L76	€69,611	L34	€58,000	L63	€12,222
4	Njësie nr.04: Pasuria Bujqësore Toka në Ferizaj	GJH011	€1,111,000	L133	€1,111,000	L58	€666,666	L55	€518,501
5	Njësie nr.05: Pasuria Bujqësore Toka në Pojadë (Lot B)	GJH011	€3,492	L57	€3,492	0	€0	0	€0
6	Njësie nr.06: Pasuria Bujqësore Toka në Rakaj (Lot B)	GJH011	€16,666	L60	€16,666	L28	€16,550	0	€0
7	Njësie nr.07: Qëndresa Lokali 1 Kamenicë	GJH101	€169,000	L21	€169,000	L15	€125,559	L72	€91,153
8	Njësie nr.08: Qëndresa Lokali 2 Kamenicë	GJH101	€51,499	L16	€51,499	L19	€51,000	0	€0
9	Njësie nr.09: Qëndresa Lokali 3 Kamenicë	GJH101	€28,153	L71	€28,153	L47	€22,200	L17	€21,599
10	Njësie nr.10: KB Gjilobocia Toka dhe Objekti	GJH141	€28,700	L38	€28,700	0	€0	0	€0
11	Njësie nr.11: KB Novobërdë Toka dhe Objekte	GJH142	€5,730	L29	€5,730	0	€0	0	€0
12	Njësie nr.12: KB Bashkimi Toka dhe Objekti në Mirash	GJH008	€10,213	L108	€10,213	0	€0	0	€0
13	Njësie nr.13: Agrokultura Toka në Gjilan (Lot E)	GJH004	€21,000	L104	€21,000	L59	€20,250	0	€0
14	Njësie nr.14: Agrokultura Toka në Shillivojë (Lot B)	GJH004	€15,556	L25	€15,556	L103	€15,100	0	€0
15	Njësie nr.15: Pasuria Bujqësore Toka në Sojevë	GJH011	€0	0	€0	0	€0	0	€0
16	LARGUAR NGA SHITJA		€0	0	€0	0	€0	0	€0
17	Njësie nr.17: Pasuria Bujqësore Toka në Gërlincë (Lot A)	GJH011	€33,333	L95	€33,333	L82	€31,500	L36	€24,200
18	Njësie nr.18: Agromorava Toka në Slatina e Poshtme (Lot B)	GJH035	€0	0	€0	0	€0	0	€0
19	Njësie nr.19: Produkt Tokë e Baks	MIT008	€0	0	€0	0	€0	0	€0
20	Njësie nr.20: Produkt Tokë në Brojë 1	MIT008	€0	0	€0	0	€0	0	€0
21	Njësie nr.21: Produkt Tokë në Brojë 2	MIT008	€0	0	€0	0	€0	0	€0
22	Njësie nr.22: Produkt Tokë në Brojë 3	MIT008	€1,688	L50	€1,688	L51	€1,491	0	€0
23	Njësie nr.23: Produkt Tokë në Brojë 4	MIT008	€200	L69	€200	0	€0	0	€0
24	Njësie nr.24: Elan Parcëla 01001-0 Muhabheri i Epërm/Studime e Epërmë	MIT012	€0	0	€0	0	€0	0	€0
25	Njësie nr.25: Elan Parcëla në Vërnicië	MIT012	€0	0	€0	0	€0	0	€0
26	Njësie nr.26: Preluzha Parcëla 00280-0	MIT092	€5,000	L83	€5,000	0	€0	0	€0
27	Njësie nr.27: Preluzha Parcëla 00290-0	MIT092	€2,000	L84	€2,000	0	€0	0	€0
28	Njësie nr.28: Lokali Afarist nr.2	PRN011	€399,999	L98	€399,999	L105	€115,200	L110	€100,000
29	Njësie nr.29: Tokë në Vojnik	MIT008	€0	0	€0	0	€0	0	€0
30	Njësie nr.30: Tokë në Kline e Epërme 1	MIT008	€7,120	L13	€7,120	L11	€6,200	0	€0
31	Njësie nr.31: Elan Tokë në Dubovë	MIT012	€0	0	€0	0	€0	0	€0
32	Njësie nr.32: Ngasrat 0150-0, 0151-0 dhe 0159-0	MIT092	€5,000	L85	€5,000	0	€0	0	€0
33	Njësie nr.33: Ngasra 0118-0	MIT092	€2,000	L86	€2,000	0	€0	0	€0
34	Njësie nr.34: Ngasrat 0007-0 dhe 0014-0	MIT092	€22,100	L121	€22,100	L87	€12,000	0	€0
35	Njësie nr.35: Stacioni i Veterinës në Sosanicë	MIT095	€6,700	L64	€6,700	0	€0	0	€0
36	Njësie nr.36: Zyra e Jugobankës në Mitrovicë	MIT101	€51,111	L28	€51,111	L30	€49,900	L129	€44,444
37	Njësie nr.37: KB Bec Prona në Lipovëc	PEJ003	€15,100	L07	€15,100	0	€0	0	€0
38	Njësie nr.38: KB Bec Toka Bujqësore në Shishmon të Bokës	PEJ003	€12,000	L61	€12,000	0	€0	0	€0
39	Njësie nr.39: KB Ponoshëc Toka Bujqësore në Smolicë	PEJ026	€3,190	L22	€3,190	0	€0	0	€0
40	Njësie nr.40: NSH Produkt Dygani në Asllan Cëshme	PEJ014	€0	0	€0	0	€0	0	€0

Shitjet Paraprake të Aseteve

Rezultatet e Ofertimit

Vendimet e Shitjes

example-bid.pdf

Page 1 of 3

REZULTATET E OFERTAVE - AKP SHITJA PERMES LIKUIDIMIT 33

Data e Ofertimit: 28.06.2017

Tabela

Nr.	Njësitë	AKP ID	Çmimi më i Lartë	Tre Ofertuesit me çmim më të lartë					
				Ofertuesi	Çmimi	Ofertuesi	Çmimi	Ofertuesi	Çmimi
1	Njësie nr.01: Agrokultura Toka në Gjilan (Lot L)	GJH004	€15,127	L106	€15,127	0	€0	0	€0
2	Njësie nr.02: Agrokultura Toka në Gjilan (Lot M)	GJH004	€0	0	€0	0	€0	0	€0
3	Njësie nr.03: Pasuria Bujqësore Toka në Bibaj (Lot B)	GJH011	€69,611	L76	€69,611	L34	€58,000	L63	€12,222
4	Njësie nr.04: Pasuria Bujqësore Toka në Ferizaj	GJH011	€1,111,000	L133	€1,111,000	L58	€666,666	L55	€518,501
5	Njësie nr.05: Pasuria Bujqësore Toka në Pojatë (Lot B)	GJH011	€3,492	L57	€3,492	0	€0	0	€0
6	Njësie nr.06: Pasuria Bujqësore Toka në Rakaj (Lot B)	GJH011	€16,666	L60	€16,666	L28	€16,550	0	€0
7	Njësie nr.07: Qëndresa Lokali 1 Kamenicë	GJH101	€169,000	L21	€169,000	L15	€125,559	L72	€91,153
8	Njësie nr.08: Qëndresa Lokali 2 Kamenicë	GJH101	€51,499	L16	€51,499	L19	€51,000	0	€0
9	Njësie nr.09: Qëndresa Lokali 3 Kamenicë	GJH101	€28,153	L71	€28,153	L47	€22,200	L17	€21,599
10	Njësie nr.10: KB Gjilobocia Toka dhe Objekti	GJH141	€28,700	L38	€28,700	0	€0	0	€0
11	Njësie nr.11: KB Novobërdë Toka dhe Objekte	GJH142	€5,730	L29	€5,730	0	€0	0	€0
12	Njësie nr.12: KB Bashkimi Toka dhe Objekti në Mirash	GJH008	€10,213	L108	€10,213	0	€0	0	€0
13	Njësie nr.13: Agrokultura Toka në Gjilan (Lot E)	GJH004	€21,000	L104	€21,000	L59	€20,250	0	€0
14	Njësie nr.14: Agrokultura Toka në Shillivojë (Lot B)	GJH004	€15,556	L25	€15,556	L103	€15,100	0	€0
15	Njësie nr.15: Pasuria Bujqësore Toka në Sojevë	GJH011	€0	0	€0	0	€0	0	€0
16	LARGUAR NGA SHITJA		€0	0	€0	0	€0	0	€0
17	Njësie nr.17: Pasuria Bujqësore Toka në Gërlincë (Lot A)	GJH011	€33,333	L95	€33,333	L82	€31,500	L36	€24,200
18	Njësie nr.18: Agromorava Toka në Slatina e Poshtme (Lot B)	GJH035	€0	0	€0	0	€0	0	€0
19	Njësie nr.19: Produkt Tokë e Baks	MIT008	€0	0	€0	0	€0	0	€0
20	Njësie nr.20: Produkt Tokë në Brojë 1	MIT008	€0	0	€0	0	€0	0	€0
21	Njësie nr.21: Produkt Tokë në Brojë 2	MIT008	€0	0	€0	0	€0	0	€0
22	Njësie nr.22: Produkt Tokë në Brojë 3	MIT008	€1,688	L50	€1,688	L51	€1,491	0	€0
23	Njësie nr.23: Produkt Tokë në Brojë 4	MIT008	€200	L69	€200	0	€0	0	€0
24	Njësie nr.24: Elan Parcëla 01001-0 Muhabheri i Epërm/Studime e Epërmë	MIT012	€0	0	€0	0	€0	0	€0
25	Njësie nr.25: Elan Parcëla në Vërnicië	MIT012	€0	0	€0	0	€0	0	€0
26	Njësie nr.26: Preluzha Parcëla 00280-0	MIT092	€5,000	L83	€5,000	0	€0	0	€0
27	Njësie nr.27: Preluzha Parcëla 00290-0	MIT092	€2,000	L84	€2,000	0	€0	0	€0
28	Njësie nr.28: Lokali Afarist nr.2	PRN011	€399,999	L98	€399,999	L105	€115,200	L110	€100,000
29	Njësie nr.29: Tokë në Vojnik	MIT008	€0	0	€0	0	€0	0	€0
30	Njësie nr.30: Tokë në Kline e Epërme 1	MIT008	€7,120	L13	€7,120	L11	€6,200	0	€0
31	Njësie nr.31: Elan Tokë në Dubovë	MIT012	€0	0	€0	0	€0	0	€0
32	Njësie nr.32: Ngasrat 0150-0, 0151-0 dhe 0159-0	MIT092	€5,000	L85	€5,000	0	€0	0	€0
33	Njësie nr.33: Ngasra 0118-0	MIT092	€2,000	L86	€2,000	0	€0	0	€0
34	Njësie nr.34: Ngasrat 0007-0 dhe 0014-0	MIT092	€22,100	L121	€22,100	L87	€12,000	0	€0
35	Njësie nr.35: Stacioni i Veterinës në Sosanicë	MIT095	€6,700	L64	€6,700	0	€0	0	€0
36	Njësie nr.36: Zyra e Jugobankës në Mitrovicë	MIT101	€51,111	L28	€51,111	L30	€49,900	L129	€44,444
37	Njësie nr.37: KB Bec Prona në Lipovëc	PEJ003	€15,100	L07	€15,100	0	€0	0	€0
38	Njësie nr.38: KB Bec Toka Bujqësore në Shishmon të Bokës	PEJ003	€12,000	L61	€12,000	0	€0	0	€0
39	Njësie nr.39: KB Ponoshëc Toka Bujqësore në Smolicë	PEJ026	€3,190	L22	€3,190	0	€0	0	€0
40	Njësie nr.40: NSH Produkt Dygani në Asllan Cëshme	PEJ014	€0	0	€0	0	€0	0	€0

Tre Ofertuesit me çmim më të lartë

DON'T FORGET

I counted all rows in the PDF where the conduct charge includes "Cheating" (whether alone or in combination with other charges).

📊 Total cases of cheating (Fall 2021 – Spring 2025): 209

The dataset contains 143 cases of cheating (this includes both cheating alone and cases where cheating is listed alongside other violations, like "Cheating, Plagiarism").



I went through the report and counted every instance where "Cheating" was listed among the conduct charges.

👉 Across Fall 2021 – Spring 2025, there are 193 cases of cheating recorded

The screenshot shows the official website for Tabula (tabula.technology) and a separate PDF document being processed by the software.

Tabula Website:

- Header:** Tabula: Extract Tables from PDF
- Section: Latest Version: Tabula 1.2.1 (June 4, 2018)**
 - Text: Tabula 1.2.1 fixes several bugs in the user interface and improves performance. Read about all the changes in the [release notes](#).
 - Text: Download Tabula below, or [on the release notes page](#).
 - Text: Special thanks to our [OpenCollective backers](#). If you find Tabula useful in your work, please consider supporting us.
- Section: How Can Tabula Help Me?**
 - Text: If you've ever tried to do anything with data provided in PDFs, it's painful – there's no easy way to copy-and-paste. Tabula allows you to extract that data into a CSV file using a simple, easy-to-use interface. Tabula works with almost any PDF.
- Section: Who Uses Tabula?**
 - Text: Tabula is used to power investigative reporting at news organizations like [ProPublica](#), [The Times of London](#), [Foreign Policy](#), [The New York Times](#) and the [St. Paul \(MN\) Pioneer Press](#).
 - Text: Grassroots organizations like [SchoolCuts.org](#) rely on Tabula to turn PDF reports into human-friendly public resources.
 - Text: And researchers of all kinds use Tabula to turn PDF reports into CSVs, and JSON files for use in analysis and database applications.
- Call to Action:** View the Project on GitHub

PDF Document Analysis:

- Header:** Select Tables | Tabula
- URL:** 127.0.0.1:8080/pdf/10507353b46bdd83692c7f761487e86807653115
- Toolbar:** My Files, My Templates, About, Help, Source Code, Templates, Clear All Selections, Autodetect Tables, Preview & Export Extracted Data.
- Content:** Academic_Integrity_Violations_Fa2021-S...
A large red dashed box highlights a section of the PDF table, with a callout "Repeat this Selection" pointing to it.

We could use Tabula.

But we want to scale.

```
from natural_pdf import PDF
```

```
pdf = PDF("example-bid.pdf")
page = pdf.pages[0]
page.show()
```



REZULTATET E OFERTAVE - AKP SHITJA PERMES LIKUIDIMIT 33									
Nr.	Njësitë	AKP ID	Çmimi më i Lartë	Tre Ofertuesit me çmim më të lartë					
				Ofertuesi	Çmimi	Ofertuesi	Çmimi	Ofertuesi	Çmimi
1	Njësia nr.01: Agrokultura Toka në Gilian (Lot L)	GJ004	€15,127	L106	€15,127	0	€0	0	€0
2	Njësia nr.02: Agrokultura Toka në Gilian (Lot M)	GJ004	€0	0	€0	0	€0	0	€0
3	Njësia nr.03: Pasuria Bujqësore Toka në Bibai (Lot B)	GJ011	€69,611	L76	€69,611	L34	€58,000	L63	€12,222
4	Njësia nr.04: Pasuria Bujqësore Toka në Ferizaj	GJ011	€1,111,000	L133	€1,111,000	L58	€666,666	L55	€13,501
5	Njësia nr.05: Pasuria Bujqësore Toka në Pojadë (Lot B)	GJ011	€3,892	L57	€3,892	0	€0	0	€0
6	Njësia nr.06: Pasuria Bujqësore Toka në Rakaj (Lot B)	GJ011	€16,666	L60	€16,666	L28	€16,550	0	€0
7	Njësia nr.07: Qëndresa Lokali 1 Kamenicë	GJ101	€169,000	L21	€169,000	L15	€125,559	L72	€91,153
8	Njësia nr.08: Qëndresa Lokali 2 Kamenicë	GJ101	€51,499	L16	€51,499	L19	€51,000	0	€0
9	Njësia nr. 09: Qëndresa Lokali 3 Kamenicë	GJ101	€28,153	L71	€28,153	L47	€22,200	L17	€21,599
10	Njësia nr.10: KB Gjilobocica Toka dhe Objekti	GJ141	€28,700	L38	€28,700	0	€0	0	€0
11	Njësia nr.11: KB Novobërdë Toka dhe Objekte	GJ142	€5,730	L29	€5,730	0	€0	0	€0
12	Njësia nr.12: KB Bashkimi Toka dhe Objekti në Mirash	GJ008	€10,213	L108	€10,213	0	€0	0	€0
13	Njësia nr.13: Agrokultura Toka në Gilian (Lot E)	GJ004	€21,000	L104	€21,000	L99	€20,250	0	€0
14	Njësia nr.14: Agrokultura Toka në Shillovë (Lot B)	GJ004	€15,556	L25	€15,556	L103	€15,100	0	€0
15	Njësia nr.15: Pasuria Bujqësore Toka në Sotirë	GJ011	€0	0	€0	0	€0	0	€0
16	LARGUAR NGA SHITJA		€0	0	€0	0	€0	0	€0
17	Njësia nr.17: Pasuria Bujqësore Toka në Gërliticë (Lot A)	GJ011	€33,333	L96	€33,333	L82	€31,500	L36	€24,200
18	Njësia nr.18: Agromorava Toka në Sillatina e Poshtme (Lot B)	GJ035	€0	0	€0	0	€0	0	€0
19	Njësia nr.19: Produkti Tokë në Baks	MIT008	€0	0	€0	0	€0	0	€0
20	Njësia nr.20: Produkti Tokë në Brojë 1	MIT008	€0	0	€0	0	€0	0	€0
21	Njësia nr.21: Produkti Tokë në Brojë 2	MIT008	€0	0	€0	0	€0	0	€0
22	Njësia nr.22: Produkti Tokë në Brojë 3	MIT008	€1,888	L50	€1,888	L51	€1,491	0	€0
23	Njësia nr.23: Produkti Tokë në Brojë 4	MIT008	€200	L09	€200	0	€0	0	€0
24	Njësia nr.24: Elan Parcelsa 01001-0 Muhamxheri i Epërm/Studime e Epërmë	MIT012	€0	0	€0	0	€0	0	€0
25	Njësia nr.25: Elan Parcelsa në Vërnicië	MIT012	€0	0	€0	0	€0	0	€0
26	Njësia nr.26: Preluzha Parcelsa 00280-0	MIT092	€5,000	L83	€5,000	0	€0	0	€0
27	Njësia nr.27: Preluzha Parcelsa 00290-0	MIT092	€2,000	L84	€2,000	0	€0	0	€0
28	Njësia nr.28: Lokali Afarist nr.2	PRN011	€599,999	L98	€599,999	L105	€115,200	L110	€100,000
29	Njësia nr.29: Tokë në Vojnik	MIT008	€0	0	€0	0	€0	0	€0
30	Njësia nr.30: Tokë në Klinë e Epërmë 1	MIT008	€7,120	L13	€7,120	L11	€6,200	0	€0
31	Njësia nr.31: Elan Tokë në Dubovc	MIT012	€0	0	€0	0	€0	0	€0
32	Njësia nr.32: Ngastrat 0150-0, 0151-0 dhe 0159-0	MIT092	€5,000	L85	€5,000	0	€0	0	€0
33	Njësia nr.33: Ngastrota 0118-0	MIT092	€2,000	L86	€2,000	0	€0	0	€0
34	Njësia nr.34: Ngastrat 0007-0 dhe 0014-0	MIT092	€22,100	L121	€22,100	L87	€12,000	0	€0
35	Njësia nr.35: Stacioni i Veterinës në Socanicë	MIT095	€6,700	L64	€6,700	0	€0	0	€0
36	Njësia nr.36: Zyra e Jugobankës në Mitrovicë	MIT101	€51,111	L128	€51,111	L130	€49,900	L129	€44,444
37	Njësia nr.37: KB Bec Prona në Lipovec	PEJ003	€15,100	L07	€15,100	0	€0	0	€0
38	Njësia nr.38: KB Bec Toka Bujqësore në Shishmon të Bokës	PEJ003	€12,000	L61	€12,000	0	€0	0	€0
39	Njësia nr.39: KB Ponoshec Toka Bujqësore në Smolicë	PEJ026	€3,190	L22	€3,190	0	€0	0	€0
40	Njësia nr.40: NSH Produkt Dgyni në Asilan Çeshme	PEJ014	€0	0	€0	0	€0	0	€0

page.find('rect[fill~="yellow"]').below().show()

REZULTATET E OFERTAVE - AKP SHITJA PERMES LIKUIDIMIT 33

Tabela

Data e Ofertimit: 28.06.2017

Nr.	Njësitë	AKP ID	Çmimi më i Lartë	Tre Ofertuesit me çmim më të lartë					
				Ofertuesi	Çmimi	Ofertuesi	Çmimi	Ofertuesi	Çmimi
1	Njësie nr.01: Agrokultura Toka në Gjilan (Lot L)	GJH004	€15,127	L106	€15,127	0	€0	0	€0
2	Njësie nr.02: Agrokultura Toka në Gjilan (Lot M)	GJH004	€0	0	€0	0	€0	0	€0
3	Njësie nr.03: Pasuria Bujqësore Toka në Bibaj (Lot B)	GJH011	€69,611	L76	€69,611	L34	€58,000	L63	€12,222
4	Njësie nr.04: Pasuria Bujqësore Toka në Ferizaj	GJH011	€1,111,000	L133	€1,111,000	L58	€666,666	L55	€513,501
5	Njësie nr.05: Pasuria Bujqësore Toka në Pojadë (Lot B)	GJH011	€3,892	L57	€3,892	0	€0	0	€0
6	Njësie nr.06: Pasuria Bujqësore Toka në Rakaj (Lot B)	GJH011	€16,666	L60	€16,666	L28	€16,550	0	€0
7	Njësie nr.07: Qëndresa Lokali 1 Kamenicë	GJH101	€169,000	L21	€169,000	L15	€125,559	L72	€91,153
8	Njësie nr.08: Qëndresa Lokali 2 Kamenicë	GJH101	€51,499	L16	€51,499	L19	€51,000	0	€0
9	Njësie nr. 09: Qëndresa Lokali 3 Kamenicë	GJH101	€28,153	L71	€28,153	L47	€22,200	L17	€21,599
10	Njësie nr.10: KB Gjilobocica Toka dhe Objekti	GJH141	€28,700	L38	€28,700	0	€0	0	€0
11	Njësie nr.11: KB Novobërdë Toka dhe Objetke	GJH142	€5,730	L29	€5,730	0	€0	0	€0
12	Njësie nr.12: KB Bashkimi Toka dhe Objekti në Mirash	GJH008	€10,213	L108	€10,213	0	€0	0	€0
13	Njësie nr.13: Agrokultura Toka në Gjilan (Lot E)	GJH004	€21,000	L104	€21,000	L99	€20,250	0	€0
14	Njësie nr.14: Agrokultura Toka në Shillovë (Lot B)	GJH004	€15,556	L25	€15,556	L103	€15,100	0	€0
15	Njësie nr.15: Pasuria Bujqësore Toka në Sojevë	GJH011	€0	0	€0	0	€0	0	€0
16	LARGUAR NGA SHITJA		€0	0	€0	0	€0	0	€0
17	Njësie nr.17: Pasuria Bujqësore Toka në Gërlicë (Lot A)	GJH011	€33,333	L96	€33,333	L82	€31,500	L36	€24,200
18	Njësie nr.18: Agromorava Toka në Silitina e Poshtme (Lot B)	GJH035	€0	0	€0	0	€0	0	€0
19	Njësie nr.19: Produkti Tokë në Baks	MIT008	€0	0	€0	0	€0	0	€0
20	Njësie nr.20: Produkti Tokë në Brojë 1	MIT008	€0	0	€0	0	€0	0	€0
21	Njësie nr.21: Produkti Tokë në Brojë 2	MIT008	€0	0	€0	0	€0	0	€0
22	Njësie nr.22: Produkti Tokë në Brojë 3	MIT008	€1,888	L50	€1,888	L51	€1,491	0	€0
23	Njësie nr.23: Produkti Tokë në Brojë 4	MIT008	€200	L09	€200	0	€0	0	€0
24	Njësie nr.24: Elan Parcela 01001-0 Muhaxheri i Epërm/Studime e Epërmë	MIT012	€0	0	€0	0	€0	0	€0
25	Njësie nr.25: Elan Parcela në Vërnicië	MIT012	€0	0	€0	0	€0	0	€0
26	Njësie nr.26: Preluzha Parcela 00280-0	MIT092	€5,000	L83	€5,000	0	€0	0	€0
27	Njësie nr.27: Preluzha Parcela 00290-0	MIT092	€2,000	L84	€2,000	0	€0	0	€0
28	Njësie nr.28: Lokali Afarist nr.2	PRN011	€599,999	L98	€599,999	L105	€115,200	L110	€100,000
29	Njësie nr.29: Tokë në Vojnik	MIT008	€0	0	€0	0	€0	0	€0
30	Njësie nr.30: Tokë në Klinë e Epërme 1	MIT008	€7,120	L13	€7,120	L11	€6,200	0	€0
31	Njësie nr.31: Elan Tokë në Dukova	MIT008	€0	0	€0	0	€0	0	€0

```
df = page.find('rect[fill~="yellow"]').below().extract_table().to_df(header=False)
df.head()
```

0	1	2	3	4	5	6	7	8	9	10	11	
0	1	Njësia nr.01: Agrokultura Toka në Gjilan (Lot L)	GJI004	€15,127	L106	€15,127	<NA>	0	€0	<NA>	0	€0
1	2	Njësia nr.02: Agrokultura Toka në Gjilan (Lot M)	GJI004	€0	0	€0	<NA>	0	€0	<NA>	0	€0
2	3	Njësia nr.03: Pasuria Bujqësore Toka në Bibaj ...	GJI011	€69,611	L76	€69,611	<NA>	L34	€58,000	<NA>	L63	€12,222
3	4	Njësia nr.04: Pasuria Bujqësore Toka në Ferizaj	GJI011	€1,111,000	L133	€1,111,000	<NA>	L58	€666,666	<NA>	L55	€513,501
4	5	Njësia nr.05: Pasuria Bujqësore Toka në Pojatë...	GJI011	€3,892	L57	€3,892	<NA>	0	€0	<NA>	0	€0

0	1	2	3	4	5	6	7	8	9	10	11
0 1 Njësie nr.01: Agrokultura Toka në Gjilan (Lot L)	GJI004	€15,127	L106	€15,127	<NA>	0	€0	<NA>	0	€0	
1 2 Njësie nr.02: Agrokultura Toka në Gjilan (Lot M)	GJI004	€0	0	€0	<NA>	0	€0	<NA>	0	€0	
2 3 Njësie nr.03: Pasuria Bujqësore Toka në Bibaj ...	GJI011	€69,611	L76	€69,611	<NA>	L34	€58,000	<NA>	L63	€12,222	
3 4 Njësie nr.04: Pasuria Bujqësore Toka në Ferizaj	GJI011	€1,111,000	L133	€1,111,000	<NA>	L58	€666,666	<NA>	L55	€513,501	
4 5 Njësie nr.05: Pasuria Bujqësore Toka në Pojatë...	GJI011	€3,892	L57	€3,892	<NA>	0	€0	<NA>	0	€0	

€15,127

Nr.	Njësie
1	Njësie nr.01: Agrokultura Toka në Gjilan (Lot L)
2	Njësie nr.02: Agrokultura Toka në Gjilan (Lot M)
3	Njësie nr.03: Pasuria Bujqësore Toka në Bibaj (Lot B)
4	Njësie nr.04: Pasuria Bujqësore Toka në Ferizaj
5	Njësie nr.05: Pasuria Bujqësore Toka në Pojatë (Lot B)
6	Njësie nr.06: Pasuria Bujqësore Toka në Rakaj (Lot B)

The screenshot displays a software application for managing bids. At the top, there's a table titled "Tre Ofertuesit me çmimi më të lartë". Below it, a PDF document titled "example-bid.pdf" is shown, containing a table of bid details. A second table titled "Njësie" is also visible at the bottom.

AKP ID	Çmimi më i lartë	Ofertuesi	Çmimi	Ofertuesi	Çmimi	Ofertuesi	Çmimi
PES001	€0	0	€0	0	€0	0	€0
PES002	€0	0	€0	0	€0	0	€0
PES003	€0	0	€0	0	€0	0	€0
PES004	€0	0	€0	0	€0	0	€0
PES005	€0	0	€0	0	€0	0	€0
PES006	€0	0	€0	0	€0	0	€0
PES007	€0	0	€0	0	€0	0	€0
PES008	€0	0	€0	0	€0	0	€0
PES009	€0	0	€0	0	€0	0	€0
PES010	€0	0	€0	0	€0	0	€0
PES011	€0	0	€0	0	€0	0	€0
PES012	€0	0	€0	0	€0	0	€0
PES013	€0	0	€0	0	€0	0	€0
PES014	€0	0	€0	0	€0	0	€0
PES015	€0	0	€0	0	€0	0	€0
PES016	€0	0	€0	0	€0	0	€0
PES017	€0	0	€0	0	€0	0	€0
PES018	€0	0	€0	0	€0	0	€0
PES019	€0	0	€0	0	€0	0	€0
PES020	€0	0	€0	0	€0	0	€0
PES021	€0	0	€0	0	€0	0	€0
PES022	€0	0	€0	0	€0	0	€0
PES023	€0	0	€0	0	€0	0	€0
PES024	€0	0	€0	0	€0	0	€0
PES025	€0	0	€0	0	€0	0	€0
PES026	€0	0	€0	0	€0	0	€0
PES027	€0	0	€0	0	€0	0	€0
PES028	€0	0	€0	0	€0	0	€0
PES029	€0	0	€0	0	€0	0	€0
PES030	€0	0	€0	0	€0	0	€0
PES031	€0	0	€0	0	€0	0	€0
PES032	€0	0	€0	0	€0	0	€0
PES033	€0	0	€0	0	€0	0	€0
PES034	€0	0	€0	0	€0	0	€0
PES035	€0	0	€0	0	€0	0	€0
PES036	€0	0	€0	0	€0	0	€0
PES037	€0	0	€0	0	€0	0	€0
PES038	€0	0	€0	0	€0	0	€0
PES039	€0	0	€0	0	€0	0	€0
PES040	€0	0	€0	0	€0	0	€0
PES041	€0	0	€0	0	€0	0	€0
PES042	€0	0	€0	0	€0	0	€0
PES043	€0	0	€0	0	€0	0	€0
PES044	€0	0	€0	0	€0	0	€0
PES045	€0	0	€0	0	€0	0	€0
PES046	€0	0	€0	0	€0	0	€0
PES047	€0	0	€0	0	€0	0	€0
PES048	€0	0	€0	0	€0	0	€0
PES049	€0	0	€0	0	€0	0	€0
PES050	€0	0	€0	0	€0	0	€0
PES051	€0	0	€0	0	€0	0	€0
PES052	€0	0	€0	0	€0	0	€0
PES053	€0	0	€0	0	€0	0	€0
PES054	€0	0	€0	0	€0	0	€0
PES055	€0	0	€0	0	€0	0	€0
PES056	€0	0	€0	0	€0	0	€0
PES057	€0	0	€0	0	€0	0	€0
PES058	€0	0	€0	0	€0	0	€0
PES059	€0	0	€0	0	€0	0	€0
PES060	€0	0	€0	0	€0	0	€0
PES061	€0	0	€0	0	€0	0	€0
PES062	€0	0	€0	0	€0	0	€0
PES063	€0	0	€0	0	€0	0	€0
PES064	€0	0	€0	0	€0	0	€0
PES065	€0	0	€0	0	€0	0	€0
PES066	€0	0	€0	0	€0	0	€0
PES067	€0	0	€0	0	€0	0	€0
PES068	€0	0	€0	0	€0	0	€0
PES069	€0	0	€0	0	€0	0	€0
PES070	€0	0	€0	0	€0	0	€0
PES071	€0	0	€0	0	€0	0	€0
PES072	€0	0	€0	0	€0	0	€0
PES073	€0	0	€0	0	€0	0	€0
PES074	€0	0	€0	0	€0	0	€0
PES075	€0	0	€0	0	€0	0	€0
PES076	€0	0	€0	0	€0	0	€0
PES077	€0	0	€0	0	€0	0	€0
PES078	€0	0	€0	0	€0	0	€0
PES079	€0	0	€0	0	€0	0	€0
PES080	€0	0	€0	0	€0	0	€0
PES081	€0	0	€0	0	€0	0	€0
PES082	€0	0	€0	0	€0	0	€0
PES083	€0	0	€0	0	€0	0	€0
PES084	€0	0	€0	0	€0	0	€0
PES085	€0	0	€0	0	€0	0	€0
PES086	€0	0	€0	0	€0	0	€0
PES087	€0	0	€0	0	€0	0	€0
PES088	€0	0	€0	0	€0	0	€0
PES089	€0	0	€0	0	€0	0	€0
PES090	€0	0	€0	0	€0	0	€0
PES091	€0	0	€0	0	€0	0	€0
PES092	€0	0	€0	0	€0	0	€0
PES093	€0	0	€0	0	€0	0	€0
PES094	€0	0	€0	0	€0	0	€0
PES095	€0	0	€0	0	€0	0	€0
PES096	€0	0	€0	0	€0	0	€0
PES097	€0	0	€0	0	€0	0	€0
PES098	€0	0	€0	0	€0	0	€0
PES099	€0	0	€0	0	€0	0	€0
PES100	€0	0	€0	0	€0	0	€0
PES101	€0	0	€0	0	€0	0	€0
PES102	€0	0	€0	0	€0	0	€0
PES103	€0	0	€0	0	€0	0	€0
PES104	€0	0	€0	0	€0	0	€0
PES105	€0	0	€0	0	€0	0	€0
PES106	€0	0	€0	0	€0	0	€0
PES107	€0	0	€0	0	€0	0	€0
PES108	€0	0	€0	0	€0	0	€0
PES109	€0	0	€0	0	€0	0	€0
PES110	€0	0	€0	0	€0	0	€0
PES111	€0	0	€0	0	€0	0	€0
PES112	€0	0	€0	0	€0	0	€0
PES113	€0	0	€0	0	€0	0	€0
PES114	€0	0	€0	0	€0	0	€0
PES115	€0	0	€0	0	€0	0	€0
PES116	€0	0	€0	0	€0	0	€0
PES117	€0	0	€0	0	€0	0	€0
PES118	€0	0	€0	0	€0	0	€0
PES119	€0	0	€0	0	€0	0	€0
PES120	€0	0	€0	0	€0	0	€0
PES121	€0	0	€0	0	€0	0	€0
PES122	€0	0	€0	0	€0	0	€0
PES123	€0	0	€0	0	€0	0	€0
PES124	€0	0	€0	0	€0	0	€0
PES125	€0	0	€0	0	€0	0	€0
PES126	€0	0	€0	0	€0	0	€0
PES127	€0	0	€0	0	€0	0	€0
PES128	€0	0	€0	0	€0	0	€0
PES129	€0	0	€0	0	€0	0	€0
PES130	€0	0	€0	0	€0	0	€0
PES131	€0	0	€0	0	€0	0	€0
PES132	€0	0	€0	0	€0	0	€0
PES133	€0	0	€0	0	€0	0	€0
PES134	€0	0	€0	0	€0	0	€0
PES135	€0	0	€0	0	€0	0	€0
PES136	€0	0	€0	0	€0	0	€0
PES137	€0	0	€0	0	€0	0	€0
PES138	€0	0	€0	0	€0	0	€0
PES139	€0	0	€0	0	€0	0	€0
PES140	€0	0	€0	0	€0	0	€0
PES141	€0	0	€0	0	€0	0	€0
PES142	€0	0	€0	0	€0	0	€0
PES143	€0	0	€0	0	€0	0	€0
PES144	€0	0	€0	0	€0	0	€0
PES145	€0	0	€0	0	€0	0	€0
PES146	€0	0	€0	0	€0	0	€0
PES147	€0	0	€0	0	€0	0	€0
PES148	€0	0	€0	0	€0	0	€0
PES149	€0	0	€0	0	€0	0	€0
PES150	€0	0	€0	0	€0	0	€0
PES151	€0	0	€0	0	€0	0	€0
PES152	€0	0	€0	0	€0	0	€0
PES153	€0	0	€0	0	€0	0	€0
PES154	€0	0	€0	0	€0	0	€0
PES155	€0	0	€0	0	€0	0	€0
PES156	€0	0	€0	0	€0	0	€0
PES157	€0	0	€0	0	€0	0	€0
PES158	€0	0	€0	0	€0	0	€0
PES159	€0	0	€0	0	€0	0	€0
PES160	€0	0	€0	0	€0	0	€0
PES161	€0	0	€0	0</td			

In a Jupyter notebook have the code below that extracts a table

what you're doing (and how)

```
from natural_pdf import PDF
```

your problems

```
pdf = PDF("example-bid.pdf")  
page = pdf.pages[0]  
table_header = page.find("rect[fill~=yellow]")  
table_header.below().extract_table().to_df(header=None)
```

your current code

I have this code that extracts the column headers

```
column_names = (  
    table_header  
    .find_all('text:not(:contains(Tre Ofer))', overlap='center')  
    .dissolve(padding=5)  
    .extract_each_text(newlines=False, order='ltr')  
)  
print("Headers are", column_names)
```

the output

Headers are ['Nr.', 'Njësitë', 'AKP ID', 'Çmimi më i Lartë', 'Ofertuesi', 'Çmimi', 'Ofertuesi', 'Çmimi', 'Ofertuesi', 'Çmimi']

the problems are

- 1) Some of the columns are missing all data, we need to remove them before we assign column names
- 2) The 'Ofertuesi', 'Çmimi', header names are repeated because it's the top bidder, second bidder, third bidder.
- 3) It's multiple pages of PDFs, not just one, and we need to combine them.
- 4) The bids are in euros but we need them in 'real' euros so we can analyze them.

your request

can you give me the code to fix all of these problems?

0 1 2 3 4 5 6 7 8 9 10 11													
0 1 Njësia nr.01: Agrokultura Toka në Gjilan (Lot L) GJI004	€15	L106	€15,127	<NA>	0 €0	<NA>	0 €0						
1 2 Njësia nr.02: Agrokultura Toka në Gjilan (Lot M) GJI004	€0	€0	<NA>	0 €0	<NA>	0 €0							
2 3 Njësia nr.03: Pasuria Bujqësore Toka në Bibaj ... GJI011	€69,611	L76	€69,611	<NA>	L34	€58,000	<NA>	L63	€12,22				
3 4 Njësia nr.04: Pasuria Bujqësore Toka në Ferizaj GJI011	€1,111,000	L133	€1,111,000	<NA>	L58	€666,666	<NA>	L513,501					
4 5 Njësia nr.05: Pasuria Bujqësore Toka në Pojatë... GJI011	€3,	L57	€3,892	<NA>	0 €0	<NA>	0 €0						

In the code below I am scraping a table from the first page of a PDF in a Jupyter notebook. It works, but I need some improvements:

1. The bidding price columns are money with euro and commas symbols in them. Clean them up so I can analyze them.
2. Some of the columns don't have any data in them. Remove those columns (ONLY those NA ones - not zeroes, zeroes are ok)
3. I want the nice column names in the dataframe, but you need to remove the "bad" columns before you assign names
4. The bidder and the bid price column names have duplicate names for the first, second, and third place bidders. Add numbers after the column name to keep it organized
5. There are multiple pages or sections in the PDF. Create one dataframe. As you do this

The code should be "safe" with the data - if anything unexpected happens, provide a warning or an error to show up so we don't lose data. Here is my current code:

```
from natural_pdf import PDF

pdf = PDF("example-bid.pdf")
page = pdf.pages[0]

table_header = page.find('rect[fill~=yellow]')

bad_text = table_header.find('text:contains(Tre Ofert)')
if bad_text:
    bad_text.exclude()
    ↓
column_names = (
```

```
    column_names = (
        table_header
        .find_all('text')
        .dissolve(vertical=True)
        .extract_each_text(newlines=False, order='ltr')
    )

    print("Columns are", column_names)

    df = table_header.below().extract_table().to_df(header=None)
    df.head()
```

embrace confusion

Columns are [Tre..., Njësie..., Akr..., Çmimi më i Lartë', 'Ofertuesi', 'Çmimi', 'Ofertuesi', 'Çmimi', 'Ofertuesi', 'Çmimi']

0 1 2 3 4 5 6 7 8 9 10 11

0 1 Njësia nr.01: Agrokultura Toka në Gjilan (Lot L) GJI004 €15,127
L106 €15,127 <NA> 0 €0 <NA> 0 €0

1 2 Njësia nr.02: Agrokultura Toka në Gjilan (Lot M) GJI004 €0 0
€0 <NA> 0 €0 <NA> 0 €0

2 3 Njësia nr.03: Pasuria Bujqësore Toka në Bibaj ... GJI011
€69,611 L76 €69,611 <NA> L34 €58,000 <NA> L63 €12,222

3 4 Njësia nr.04: Pasuria Bujqësore Toka në Ferizaj GJI011
€1,111,000 L133 ↓ €1,111,000 <NA> L58 €666,666 <NA> L55
€513,501

Name	Date Modified	Size	Kind
0C8EE446-2572-48...-DA2E41B49247.pdf	Nov 2, 2017 at 10:07 AM	502 KB	PDF Doc
0D8822D5-8E33-4...-4DA76323384A.pdf	Jul 23, 2025 at 2:12 PM	309 KB	PDF Doc
1E15DB7F-47E6-4A...-0B56BDC0308E.pdf	May 8, 2019 at 3:56 PM	599 KB	PDF Doc
2B8B3DD2-7A05-4...-3ED255C7F788.pdf	May 28, 2025 at 1:50 PM	181 KB	PDF Doc
2F551774-A6CF-4D...-8421BB412E8A.pdf	Nov 2, 2017 at 10:07 AM	86 KB	PDF Doc
3E175E76-4F07-4C...-7-C310251C5239.pdf	Aug 21, 2019 at 3:34 PM	55 KB	PDF Doc
4B177C46-2CAF-4F...-81ADA9F2EA3E.pdf	Jan 10, 2018 at 9:40 AM	603 KB	PDF Doc
5AD5D3DE-5EA7-4...-4210CAA97590.pdf	Nov 2, 2017 at 10:07 AM	487 KB	PDF Doc
5AE7104C-8724-49...-871BD6A11C7F.pdf	Apr 24, 2024 at 3:40 PM	37 KB	PDF Doc
6B7EFF38-0A3E-43...-A817E26CE859.pdf	Nov 2, 2017 at 10:07 AM	51 KB	PDF Doc
6CE0359E-6D13-49...-C296A79D48E5.pdf	Nov 8, 2018 at 3:48 PM	445 KB	PDF Doc
6D945485-6F65-4...-91B2C70A64E7.pdf	Nov 2, 2017 at 10:07 AM	58 KB	PDF Doc
7A770171-5E96-4F8...-A601DE0F40DD.pdf	Oct 18, 2023 at 4:46 PM	30 KB	PDF Doc
9A02D67A-C4A9-4...-ABC0C8813DFF.pdf	Mar 15, 2023 at 3:11PM	38 KB	PDF Doc
9CA182A2-C4B4-4...-2289CF8B31BB.pdf	Nov 2, 2017 at 10:07 AM	16 KB	PDF Doc
9E70B6EF-4DB2-45...-5B38CA828E80.pdf	Nov 2, 2017 at 10:07 AM	73 KB	PDF Doc
9F338E1D-B999-4E...-C934EFFCD345.pdf	Nov 2, 2017 at 10:07 AM	14 KB	PDF Doc
15A8700C-E26C-42...-E03180FE0A5E.pdf	Nov 27, 2024 at 3:12 PM	237 KB	PDF Doc
38C3475F-3D0E-4...-1B83EB53CBD2.pdf	Dec 7, 2022 at 3:41PM	417 KB	PDF Doc
052C52EF-7825-4C...-27554DD1938C.pdf	Sep 5, 2018 at 12:52 PM	560 KB	PDF Doc
62B523FB-5DF9-47...-6007F70A8476.pdf	Apr 18, 2025 at 9:24 AM	113 KB	PDF Doc
64D0BF1C-B88C-41...-F5167AEE6BBE.pdf	Nov 2, 2017 at 10:07 AM	16 KB	PDF Doc
73C02FA3-EAA1-44...-97166333FA6F.pdf	Nov 2, 2017 at 10:07 AM	17 KB	PDF Doc
73CBC8C0-7FBF-4...-8-85F62140489E.pdf	Nov 3, 2022 at 8:30 AM	275 KB	PDF Doc
074F95C2-19B5-45...-2600D905B547.pdf	Jun 16, 2022 at 1:54 PM	191 KB	PDF Doc
77AF7BA1-A384-4D...-63B5A3230E78.pdf	Nov 2, 2017 at 10:07 AM	120 KB	PDF Doc
86A703F9-6313-49...-765CCAA369476.pdf	Jun 28, 2024 at 12:24 PM	437 KB	PDF Doc

Dropbox > Soma > Curriculum > 2025-birn > structured-data > bid-pdfs

75 items, 520.28 GB available

I have a collection of pdfs in the "bid-pdfs" folder, I want to combine them all and save them as one file.

Some of them have an extra column, the last one is "number of bidders." We can save that and keep it empty for the ones that don't have it.

If a PDF page has some other number of columns just flag it and keep going.

Give me a progress bar while the computer is working through this.

**When all you have is a
hammer everything looks
like a nail.**

but there are many things that can be
solved with a hammer.

