

Exposing the Dark Side of Google Dorks: How I Extracted Millions of Data.



Pwndec0c0

Follow

Published in
OSINT Team

4 min read

Nov 11, 2024

Listen

Share

More

How I was able to scrape large data of emails from the top mail domains(gmail,yahoo,hotmail etc..), .gov, .edu on of all countries and with private email as well with just using 2 simple tools **google dorks** and **bash** Let's start with the how can this be useful... Beware of my words I'll be using borrow,accidentally, not intended trust me its true 😊 Email scraping, is useful for legitimate purposes like marketing campaigns and data analysis, can also be exploited for malicious activities such as spam,phishing attacks, malware distribution.



I'll show you how easily you can gather millions of emails just using google dork/ google hacking. Now what's google dork?? google it! don't be lazy but if you're already aware of it Noiceeeee. here's a sample of the google dork I used.

intext:@"yahoo|gmail|hotmail".com filetype:txt site:.us

Google

intext:@"yahoo|gmail|hotmail".com filetype:txt site:.us

All Images Videos News Shopping Web Books More Tools

office365.us https://www.office365.us/app-ads/ app-ads.txt

... com with any questions #use Microsoft Ads.txt Intake Form v1 when requesting changes to yahoo.com, 57872, RESELLER #inmobi outrbrain.com ...

Los Alamos County https://www.lasalamosnm.us/files/recruitment-programs/ 2024-vendor-meeting-chat.txt

... gmail.com 01:07:09 Debbie Debbie DeSimone Optimal Performance s@protonmail.com 01:09:01 Ana Ana Martinez - Food Vendor...

dot.state.tx.us https://ftp.dot.state.tx.us/bdot-info/cmd/cserve/b... bqlist.txt

... YAHOO.COM 18230 A & G LANDSCAPE AND IRRIGATION LLC: 75203 (214) 406-2557 20@GMAIL... HOTMAIL.COM 16960 AZA ...

sos.state.mn.us https://electionresultsfiles.sos.state.mn.us/CandTbl/ CandTbl.txt

... gmail.com 01020201; Oliver Steinberg 0102; U.S. Senator; 88; 02; GLC; 1503 ... yahoo.com 01020404; Paula Overby 0102; U.S. Senator; 88; 04; DFL; NOT REQUIRED ...

Colorado County https://www.co.colorado.tx.us/upload/page/ Work Schedule - Post.txt

... hotmail.com Alt Judge Billy Kahn 651-428-9800 " JS, TX 78893" @yahoo.com Clerk n 979-942-0301 "11 IE ...

Results will give you sensitive information already
Now we can manually check and grab those names, emails,

addresses and other sensitive information for the good thing we're planning to use it



But come on!? “Manually??” we’re going to automate this using bash for our sample; We’ll be using **Google’s custom search** this will require for you to create an account and generate your own API Keys. *But we’re not going to do that* we’re just going to “borrow” someones API key for the sake of showing this sample, I’ll be dorking this as well... and let’s see if we can find one

```
① https://www.googleapis.com › customsearch › v1?key=Alza... 14IW... ...
www.googleapis.com
{
  "kind": "customsearch#search",
  "url": {
    "type": "application/json",
    "template": "https://www.googleapis.com/customsearch/v1?q=[searchTerms]&num=[count?]&start ...

```

Ohhh look I accidentally found one!
now we’re going to test this via curl to
curl "https://www.googleapis.com/customsearch/v1?
key=[API KEY]&cx=[CX]u&q=[SEARCH STRING]"

```
pwndec0c0@MarcU-LTP: ~ + | x
pwndec0c0@MarcU-LTP:~$ curl "https://www.googleapis.com/customsearch/v1?key=AIza... &cx=0
  i&q=pwndec0c0"
{
  "kind": "customsearch#search",
  "url": {
    "type": "application/json",
    "template": "https://www.googleapis.com/customsearch/v1?q={searchTerms}&num={count}&start={startIndex}&lr={language?}&safe={safe?}&cx={cx?}&sort={sort?}&filters={filter?}&gl={gl?}&cr={cr?}&googleHost={googleHost?}&c2coff={disableCnTwTranslation?}&hq={hq?}&hl={hl?}&siteSearch={siteSearch?}&siteSearchFilter={siteSearchFilter?}&exactTerms={exactTerms?}&excludeTerms={excludeTerms?}&linkSite={linkSite?}&orTerms={orTerms?}&dateRestrict={dateRestrict?}&lowRange={lowRange?}&highRange={highRange?}&searchType={searchType?}&fileType={fileType?}&rights={rights?}&imgSize={imgSize?}&imgType={imgType?}&imgColorType={imgColorType?}&imgDominantColor={imgDominantColor?}&alt=json"
  },
  "queries": {
    "request": [
      {
        "title": "Google Custom Search - pwndec0c0",
        "totalResults": "8820",
        "searchTerms": "pwndec0c0",
        "count": 10,
        "startIndex": 1,
        "inputEncoding": "utf8",
        "outputEncoding": "utf8",
        "safe": "off",
        "cx": "0"
      }
    ],
    "nextPage": [
      {
        "title": "Google Custom Search - pwndec0c0",
        "totalResults": "8820",
        "searchTerms": "pwndec0c0",
        "count": 10,
        "startIndex": 11,
        "inputEncoding": "utf8",
        "outputEncoding": "utf8",
        "safe": "off",
        "cx": "0"
      }
    ]
  }
}
```

Ohhhh it worked!

Now we're going to create our script to gather emails or any data from google.

#Search String

intext:@"yahoo|gmail|hotmail".com filetype:txt site:.us

#we are checking for any webpages for the pattern above on google##!/bin/bash

#we are getting the number of google pages results here

```
PAGECNT=$(curl -s "https://www.googleapis.com/
customsearch/v1?q=[SEARCH STRING]&key=[API
KEY]&cx=[cx]&start=1" | jq -r '.queries.nextPage[0].startIndex')
#display total number
echo "$PAGECOUNT"
```

#we're going to iterate in those pages to gather all URL results and save it

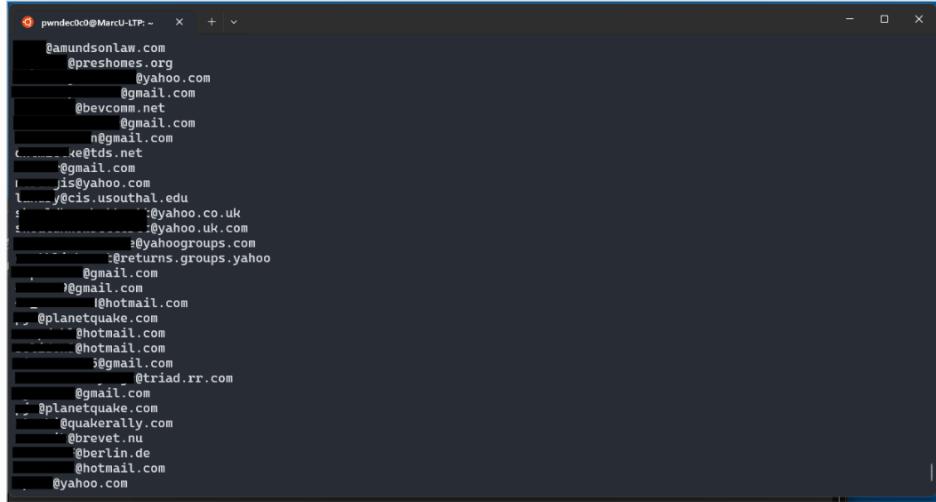
for i in (seq 1 \$PAGECNT)

do

```
curl -s 'https://www.googleapis.com/customsearch/v1?
q=[SEARCH STRING]&num=10&key=[API
KEY]&cx=[CX]&start='$i | jq -r '.items[].link' | anew links.txt
```

done

```
#we're going to iterate on those URLs and check for email  
patterns using regex  
for link in `cat links.txt`;do  
curl -s $link |grep -E -o '[a-zA-Z0-9._%+-]+@[a-zA-Z0-9.-]+\.[a-zA-Z]{2,}' | tee -a emails.txt  
done
```



The terminal window shows a list of email addresses extracted from various websites. The list includes:

- @amundsonlaw.com
- @preshomes.org
- @yahoo.com
- @gmail.com
- @bevcomm.net
- @gmail.com
- n@gmail.com
- ke@tds.net
- @gmail.com
- jis@yahoo.com
- ycis.usouthal.edu
- @yahoo.co.uk
- @yahoo.uk.com
- @yahoomgroups.com
- @returns.groups.yahoo
- @gmail.com
- @hotmail.com
- @planetquake.com
- @hotmail.com
- @hotmail.com
- i@gmail.com
- @triad.rr.com
- @gmail.com
- @planetquake.com
- @quakerally.com
- @brevet.nu
- @berlin.de
- @hotmail.com
- @yahoo.com

Guess what!? we're already able to fetch these emails

```
pwndec0c0@MarcU-LTP:~$ cat emails.txt | wc -l | sort -u  
202378
```

we were able to gather 200k email with just using google dorks and bash

Now as you can see we're just limiting our test on text files we can expand this to other filetypes as well including csv,log files, back up files etc this will give you more data, and since the site focuses on ".us" only this will isolate our scrapping of emails with only to domains that contains ".us" we can also change this to expand our data gathering you can also set this for specific target domains to check for any leaks.

As a Penetration Tester this is a great tool to have specially if your doing blackbox testing you can isolate to the targets domain and bruteforce those accounts for any weak passwords. You can improve the script that I showed you

fine tune it, run it on thread, check for emails validity etc.. customize it to help you on your use case Just remember do it ethically ;).

Take Aways

alway use “**Plus addressing**” when using your email to register, subscribe or anything related to you giving your email. E.g subscribing to Netflix use “youremail1337+netflix.com” no worries you can still receive your email in that way, advantage of this is once the email you received that was intended to “youremail1337+netflix.com” does not come from Netflix you are already aware that your email from Netflix was leaked or sold to third parties(I’m not saying Netflix sells your data! come on guys! or is it? LOL) anyways hope this help :) Chow! Thanks,

————— ✨ Support me ✨ ———

Follow me on X : [https://x.com/pwendec0c0](https://x.com/pwndec0c0)

Subscribe to my YT : <https://www.youtube.com/@pwendec0c0Tv>

Buy me a coffee :