



# Data-Driven Threat Intelligence: Useful Methods and Measurements for Handling Indicators (#ddti)

Alex Pinto  
Chief Data Scientist  
Niddel / MLSec Project  
@alexcpsc  
@MLSecProject

Alexandre Sieira  
CTO  
Niddel  
@AlexandreSieira  
@NiddelCorp

# MLSec Project / Niddel

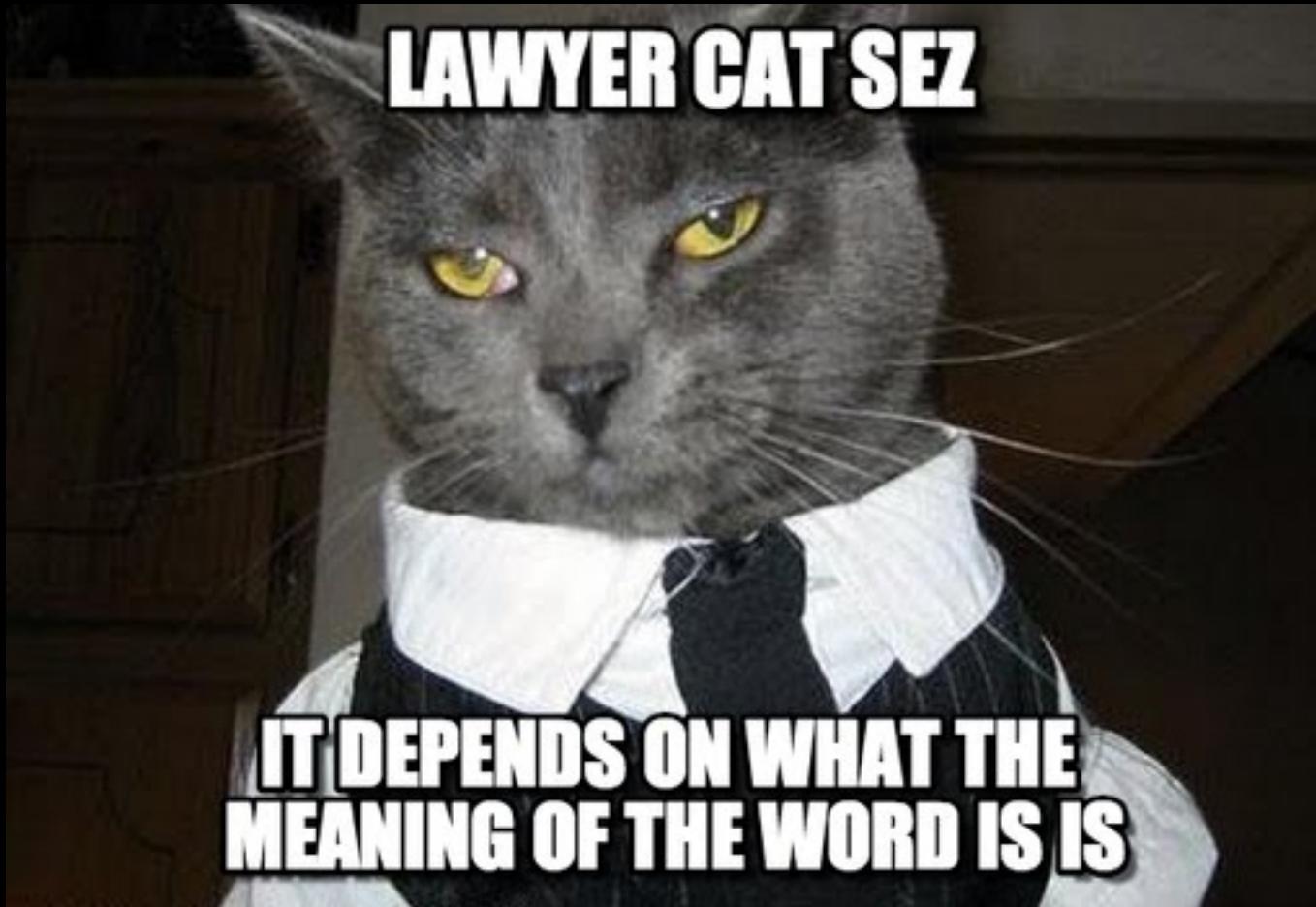
- MLSEC Project – research-focused branch of Niddel for open-source tools and community building
- Niddel builds Magnet, an Applied Threat Intelligence Platform, that leverages ML models to facilitate the usage of TI in a fully personalized and hassle-free way.
- Looking for beta testers and research collaboration
- More info at:
  - [niddel.com](http://niddel.com)
  - [mlsecproject.org](http://mlsecproject.org)



# Agenda

- ~~Cyber War...~~ Threat Intel – What is it good for?
- Combine and TIQ-test
- Using TIQ-test
  - Novelty Test
  - Overlap Test
  - Population Test
  - Aging Test
  - Uniqueness Test
- Use case: Feed Comparison

# What is TI good for anyway?



# What is TI good for anyway?

- 1) Attribution



# What is TI good for anyway?

The screenshot shows a web browser window with the URL [sony.attributed.to](http://sony.attributed.to). The page content is as follows:

**TLP: White**

## Sony breach linked to Romanian external activist group

### Executive Summary

On November 24, 2014, personally identifiable information about Sony Pictures Entertainment (SPE) employees and their dependents, e-mails between employees, information about executive salaries at the company, copies of unreleased Sony films, and other information, was obtained and released by a hacker group going under the moniker "Guardians of Peace" or "GOP".

Although the motives for the hack have yet to be revealed, the hack has been tied to the planned release of the film *The Interview*, which depicts an assassination attempt on North Korean leader Kim Jong-un, with the hackers threatening acts of terrorism if the film were to be released.

Recently, a team of 2 researchers from iDefense examined the evidence left behind by the attackers. This research has provided insight into the likely source of these attacks. Though not definitive, our analysis provides a much clearer picture and suggests an external activist group operating out of Romania is responsible for the data breach impacting Sony Pictures Entertainment. This disclosure casts further doubt on the FBI's assertion that the attack was carried out by state-sponsored actors under the control of North Korea, a theory that has been all but discredited by a host of security professionals since the attack became public, including security product pre-sales engineer Nellie Nau.

Our product indicates a different, more sinister source behind the Sony attack.  
— Nellie Nau, security product pre-sales engineer

The research team is quite certain, however, that the Guardians of Peace hacker group played no role in this attack. The clues left behind confirm that the group claiming responsibility were a fabrication to throw investigators off the trail and to mask the true source.

### Links to Romania

The research team was able to reconstruct the attack from the ground up and discovered a number of IP addresses that are linked to other attacks that have been attributed to actors in Romania as well as the presence of Romanian text in the comment strings of the malware that was recovered during the forensic investigation. Some of these malware samples have also been used in Romanian attacks.

Additional signals intelligence acquired by the research team has also implicated an actor based in Romania. This intelligence is highly classified and cannot be released in a public document, but the research team has briefed investigators with the U.S. Federal Bureau of Investigation on their findings.

# What is TI good for anyway?

- 2) Cyber Threat Maps

(<https://github.com/hrbrmstr/pewpew>)



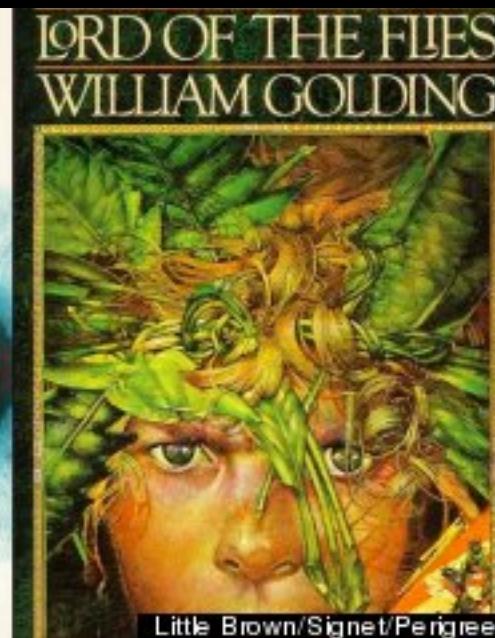
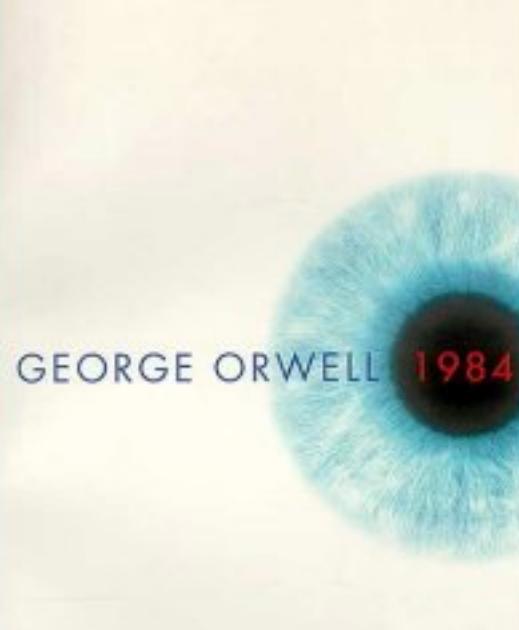
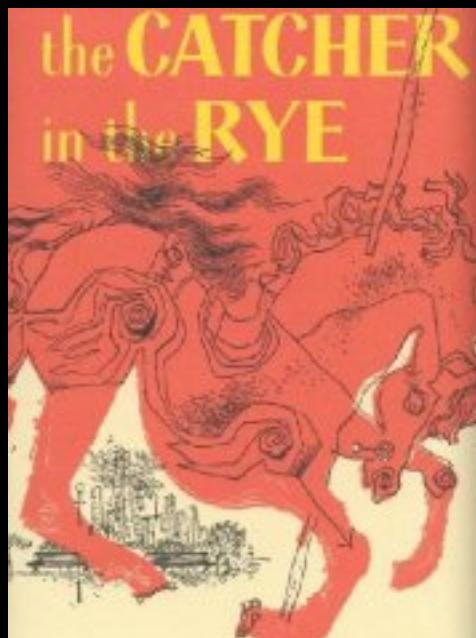
# What is TI good for anyway?

- 3) How about actual defense?
  - Use it as blacklists? As research data?
  - Thing is RAW DATA is hard to work with



# (Semi-)Required Reading

- #tiqtest slides: <http://bit.ly/tiqtest>
- RPubs Page: <http://bit.ly/tiqtest-rpubs>



# Combine and TIQ-Test

- Combine (<https://github.com/mlsecproject/combine>)
  - Gathers TI data (ip/host) from Internet and local files
  - Normalizes the data and enriches it (AS / Geo / pDNS)
  - Can export to CSV, “tiq-test format” and CRITs
  - Coming soon: CybOX / STIX (ty @kylemaxwell)
- TIQ-Test (<https://github.com/mlsecproject/tiq-test>)
  - Runs statistical summaries and tests on TI feeds
  - Generates charts based on the tests and summaries
  - Written in R (because you should learn a stat language)

# Using TIQ-TEST

- Available tests and statistics:
  - NOVELTY – How often do they update themselves?
  - OVERLAP – How do they compare to what you got?
  - POPULATION – How does this population distribution compare to another one ?
  - AGING – How long does an indicator sit on a feed?
  - UNIQUENESS – How many indicators are found in only one feed?

# Using TIQ-TEST

- New dataset!
- <https://github.com/mlsecproject/tiq-test-Winter2015>

```
print(tiq.data.getAvailableDates("raw", "public_outbound"))
```

```
## [1] "20141001" "20141002" "20141003" "20141004" "20141005" "20141006"
## [7] "20141007" "20141008" "20141009" "20141010" "20141011" "20141012"
## [13] "20141013" "20141014" "20141015" "20141016" "20141017" "20141018"
## [19] "20141019" "20141020" "20141021" "20141022" "20141023" "20141024"
## [25] "20141025" "20141026" "20141027" "20141028" "20141029" "20141030"
## [31] "20141031" "20141101" "20141102" "20141103" "20141104" "20141105"
## [37] "20141106" "20141107" "20141108" "20141109" "20141110" "20141111"
## [43] "20141112" "20141113" "20141114" "20141115" "20141116" "20141117"
## [49] "20141118" "20141119" "20141120" "20141121" "20141122" "20141123"
## [55] "20141124" "20141125" "20141126" "20141127" "20141128" "20141129"
## [61] "20141130"
```

# Using TIQ-TEST – Feeds Selected

- Dataset was separated into “inbound” and “outbound”

```
inbound.ti = tiq.data.loadTI("raw", "public_inbound", "20141101")
unique(inbound.ti$source)
```

```
## [1] "alienvault"          "autoshun"           "blocklistde"
## [4] "botscout"             "bruteforceblocker" "charleshaley"
## [7] "ciarmy"               "dragonresearch"    "dshield"
## [10] "honeypot"              "openbl"             "packetmail"
## [13] "virbl"
```

```
outbound.ti = tiq.data.loadTI("raw", "public_outbound", "20141101")
unique(outbound.ti$source)
```

```
## [1] "alienvault"          "feodo"              "malcode"
## [4] "malcode_zones"         "malwaredomainlist" "malwaredomains"
## [7] "malwaregroup"          "palevotracker"      "spyeye"
## [10] "sslbl"                "zeus"
```

# Using TIQ-TEST – Data Prep

- Extract the “raw” information from indicator feeds
- Both IP addresses and hostnames were extracted

```
outbound.ti = tiq.data.loadTI("raw", "public_outbound", "20141101")
outbound.ti[, list(entity, type, direction, source, date)]
```

```
##                                     entity type direction      source      date
## 1:          1.168.15.140  IPv4  outbound alienvault 2014-11-01
## 2:          1.93.6.86   IPv4  outbound alienvault 2014-11-01
## 3:          100.42.211.4  IPv4  outbound alienvault 2014-11-01
## 4:          101.227.172.24  IPv4  outbound alienvault 2014-11-01
## 5:          101.36.81.55  IPv4  outbound alienvault 2014-11-01
## ----
## 11388:      up.frigobar.it FQDN  outbound      zeus 2014-11-01
## 11389:      update.odeen.eu FQDN  outbound      zeus 2014-11-01
## 11390: update.rifugiopestese.it FQDN  outbound      zeus 2014-11-01
## 11391: vahendkarasis4.com FQDN  outbound      zeus 2014-11-01
## 11392: welcahlllyn.com FQDN  outbound      zeus 2014-11-01
```

# Using TIQ-TEST – Data Prep

- Convert the hostname data to IP addresses:
  - Active IP addresses for the respective date (“A” query)
  - Passive DNS from Farsight Security (DNSDB)
- For each IP record (including the ones from hostnames):
  - Add asnumber and asname (from MaxMind ASN DB)
  - Add country (from MaxMind GeoLite DB)
  - Add rhost (again from DNSDB) – most popular “PTR”

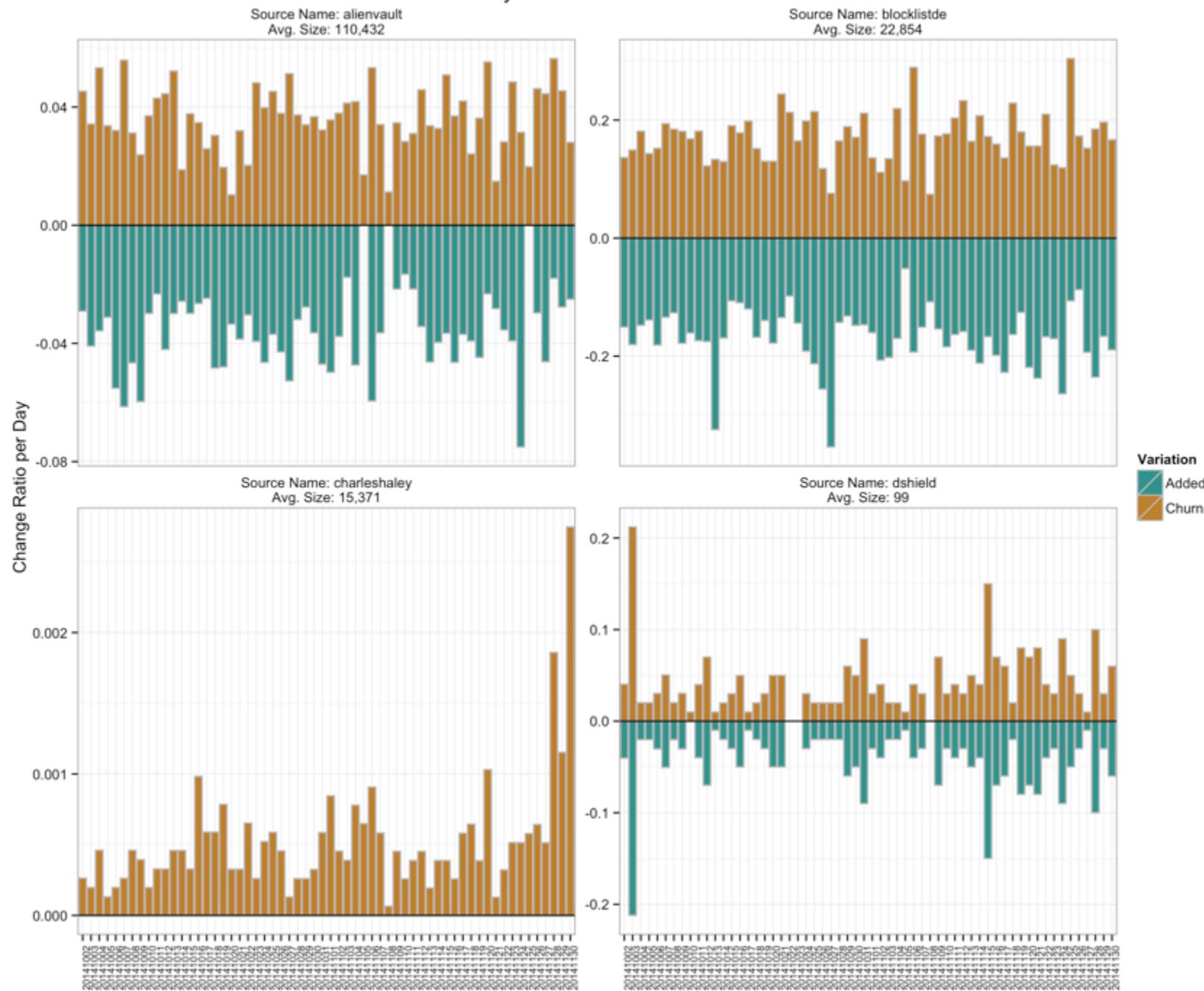
# Using TIQ-TEST – Data Prep Done

```
enrich.ti = tiq.data.loadTI("enriched", "public_outbound", "20141101")
enrich.ti = enrich.ti[, notes := NULL]
tail(enrich.ti)
```

```
##           entity type direction source      date asnumber
## 1: 94.102.63.153 IPv4 outbound zeus 2014-11-01    29073
## 2: 94.103.36.55 IPv4 outbound zeus 2014-11-01    47894
## 3: 95.163.121.12 IPv4 outbound zeus 2014-11-01   12695
## 4: 98.131.185.136 IPv4 outbound zeus 2014-11-01   32392
## 5: 98.131.185.136 IPv4 outbound zeus 2014-11-01   32392
## 6: 99.181.5.83 IPv4 outbound zeus 2014-11-01    7018
##           asname country          host
## 1: Ecatel Network     NL            NA
## 2: VeriTeknik Bilişim Ltd. TR            NA
## 3: Digital Networks CJSC    RU            NA
## 4: Ecommerce Corporation US            NA
## 5: Ecommerce Corporation US projects.globaltronics.net
## 6: AT&T Services, Inc.  US            NA
##           rhost
## 1: exadomains.net
## 2: datacenter.veriteknik.com
## 3: NA
## 4: NA
## 5: NA
## 6: adsl-99-181-5-83.dsl.irvnca.sbcglobal.net
```

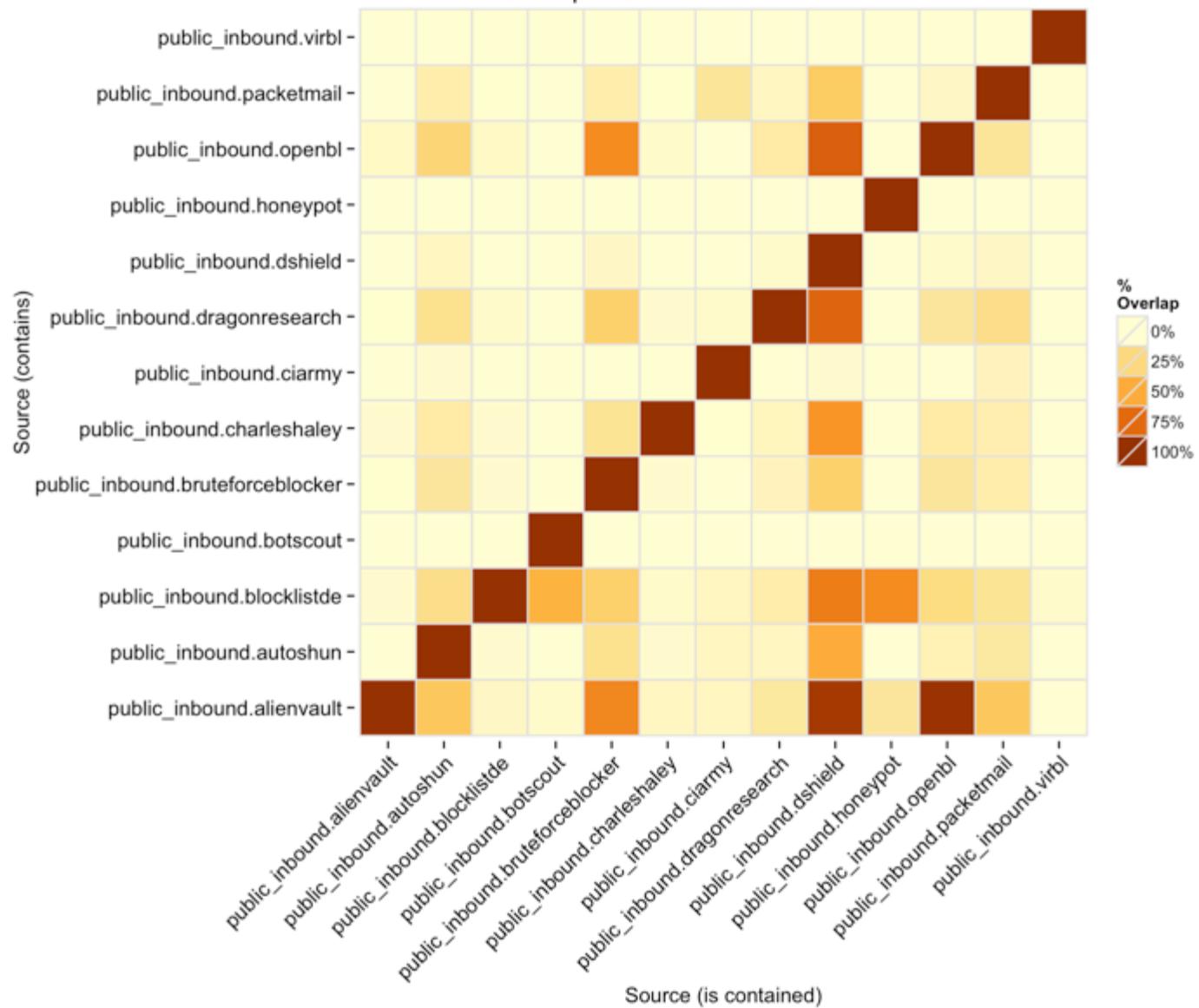
*Novelty Test* – measuring added  
and dropped indicators

### Novelty Test - Inbound Indicators

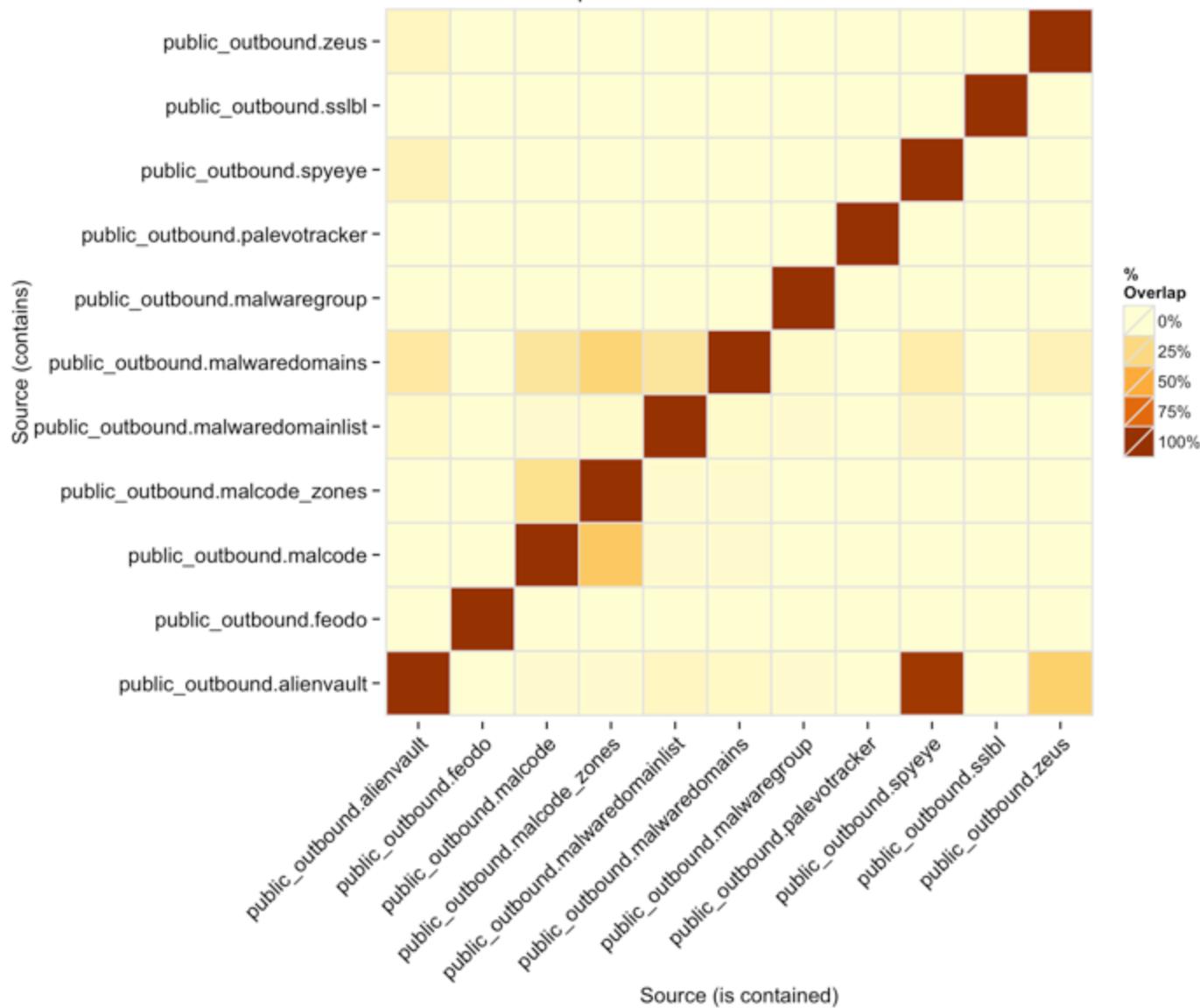


Overlap Test – More data is better,  
but make sure it is not the same  
data

Overlap Test - Inbound Data - 20141101

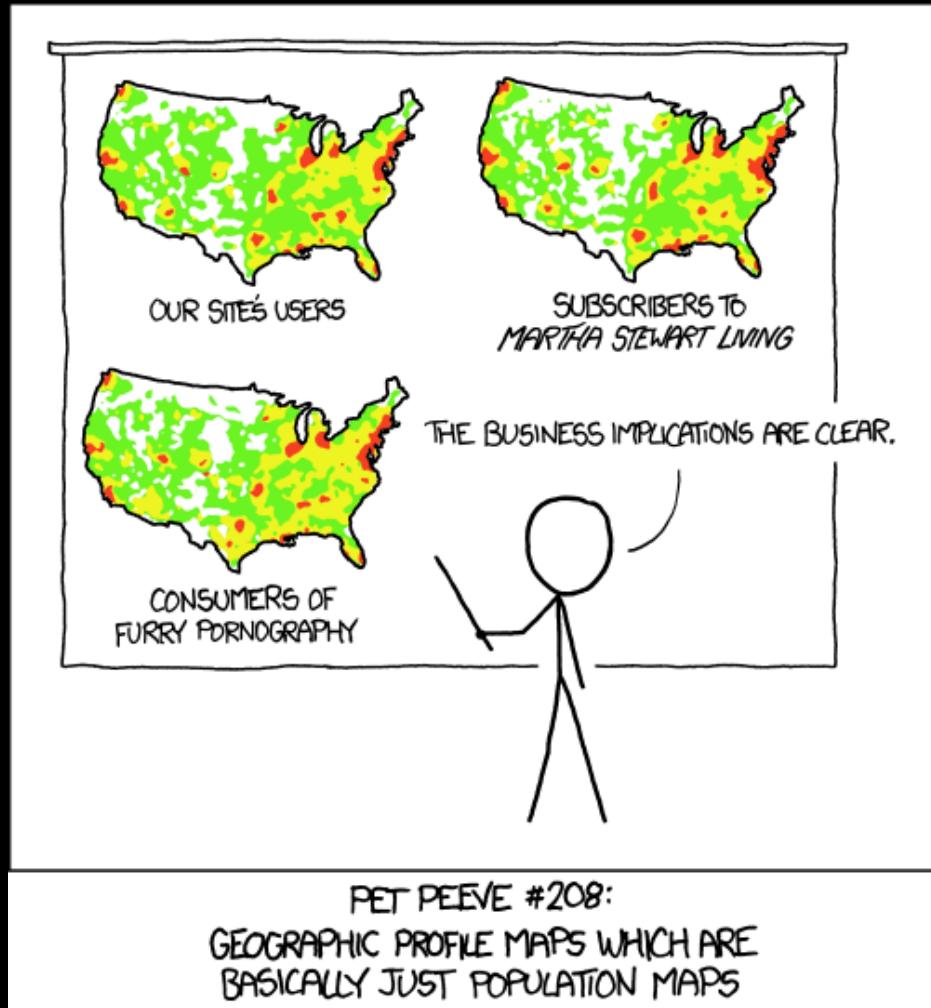


Overlap Test - Outbound Data - 20141101

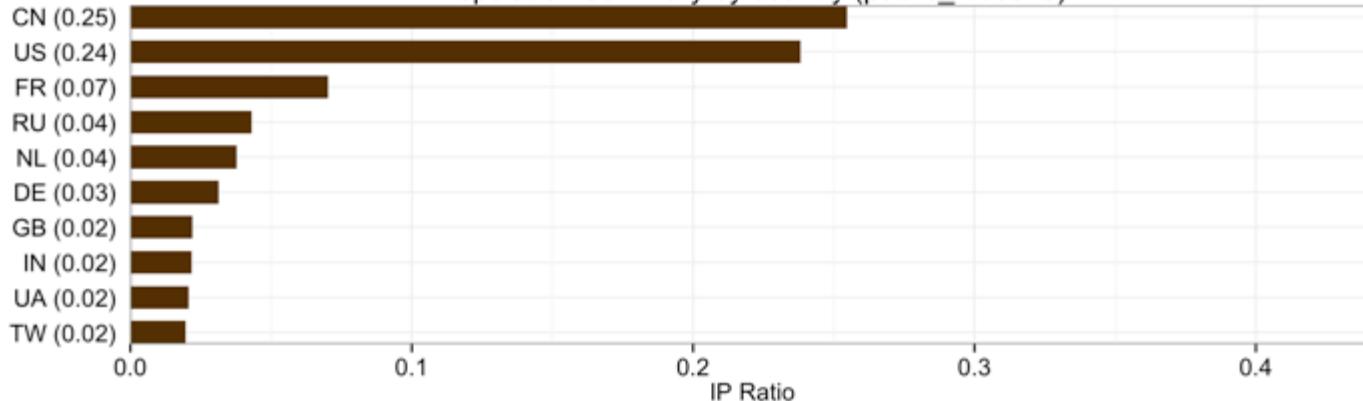


# Population Test

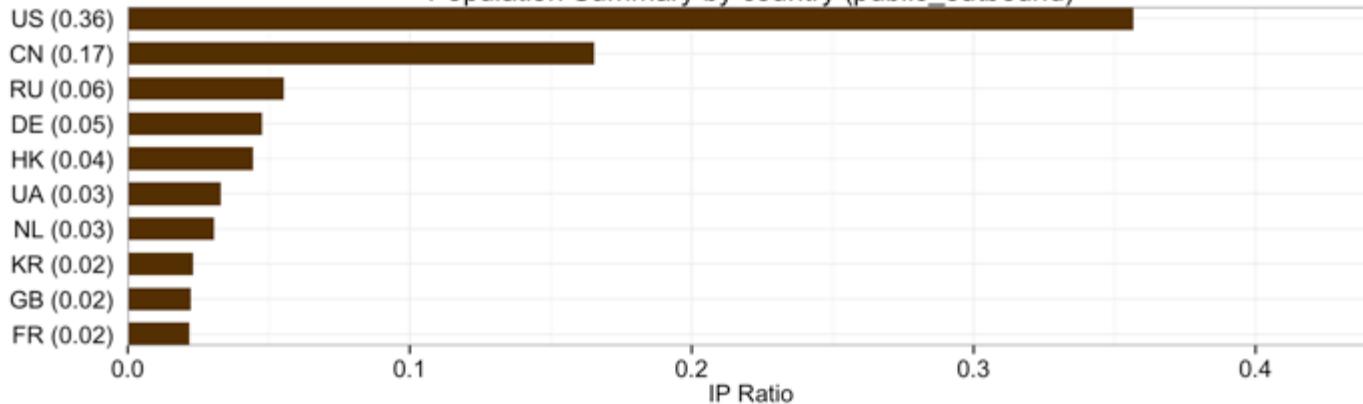
- Let us use the ASN and GeolP databases that we used to enrich our data as a reference of the “true” population.
- But, but, human beings are unpredictable! We will never be able to forecast this!



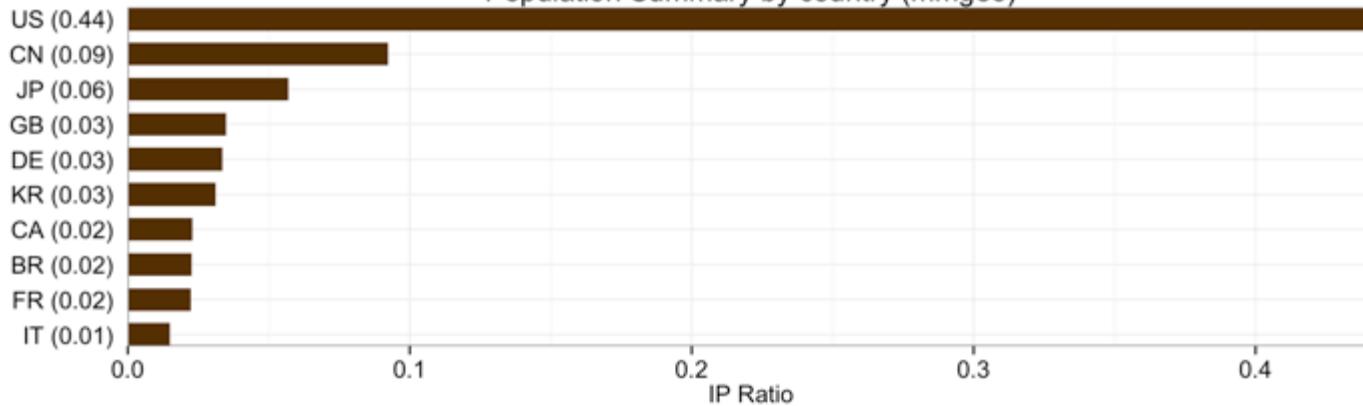
Population Summary by country (public\_inbound)



Population Summary by country (public\_outbound)



Population Summary by country (mmgeo)



# Is your sampling poll as random as you think?



# Can we get a better look?

- Statistical inference-based comparison models (hypothesis testing)
  - Exact binomial tests (when we have the “true” pop)
  - Chi-squared proportion tests (similar to independence tests)

$$(\sqrt{(-things)})^2$$

**THINGS JUST GOT REAL.**

```
tests = tiq.test.populationInference(complete.pop$mmgeo,
                                      outbound.pop$public_outbound, "country",
                                      exact = TRUE, top=10)
```

*# Whose proportion is bigger than it should be?*

```
##   country conf.int.start conf.int.end    p.value
## 1:     TH      0.047044      0.05415  0.000e+00
## 2:     US      0.025335      0.04111  9.406e-17
## 5:     HK      0.014238      0.01868 2.412e-128
## 6:     NL      0.007818      0.01268  4.091e-23
```

*# Whose is smaller?*

```
tests[p.value < 0.05/10 & conf.int.start < 0][order(conf.int.start, decreasing=F)]
```

```
##   country conf.int.start conf.int.end    p.value
## 1:     GB      -0.01926     -0.015040 8.988e-38
## 2:     CN      -0.01469     -0.005996 5.356e-06
```

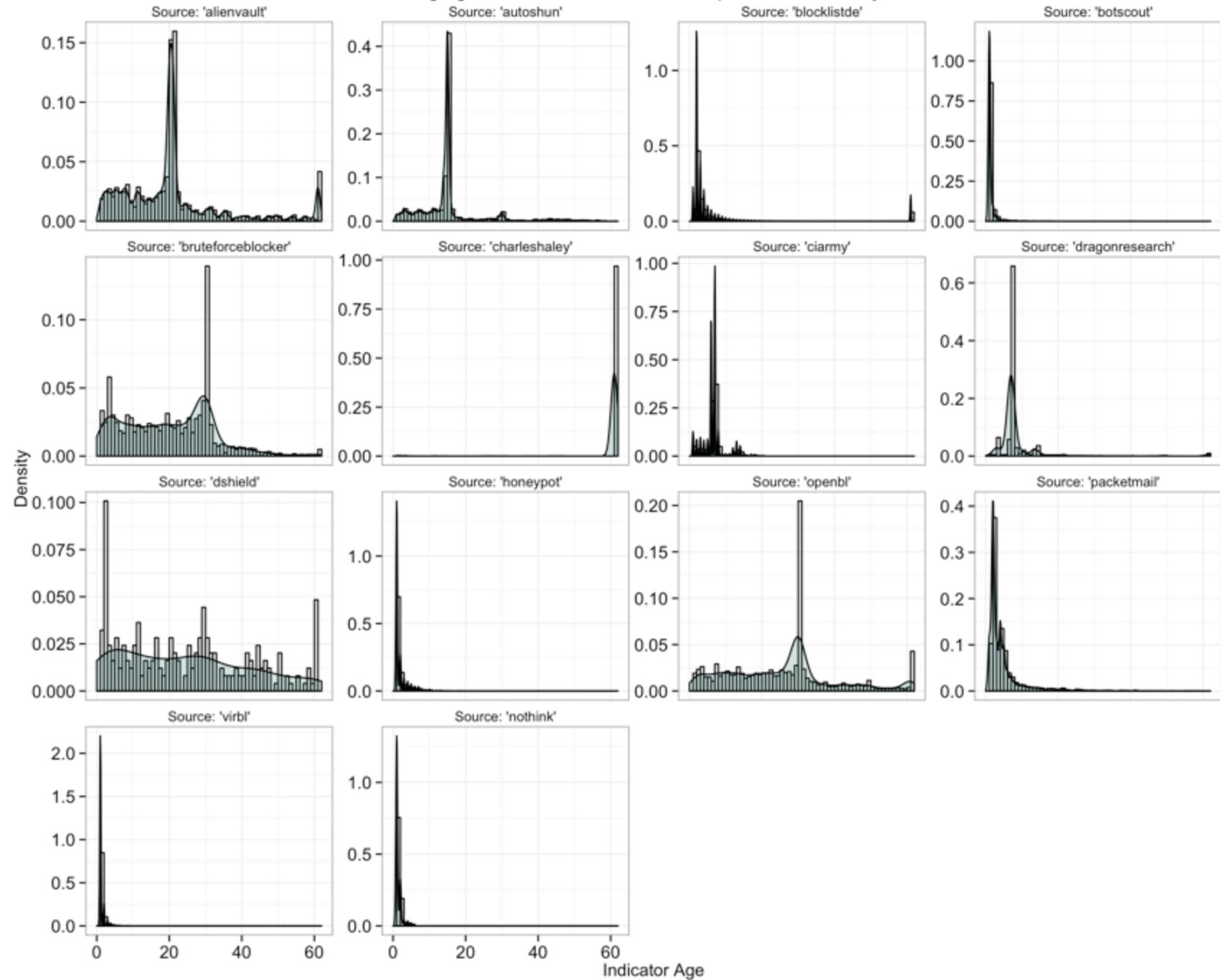
*# And whose is the same? ^\_^(ツ)\_/~*

```
tests[p.value > 0.05/10]
```

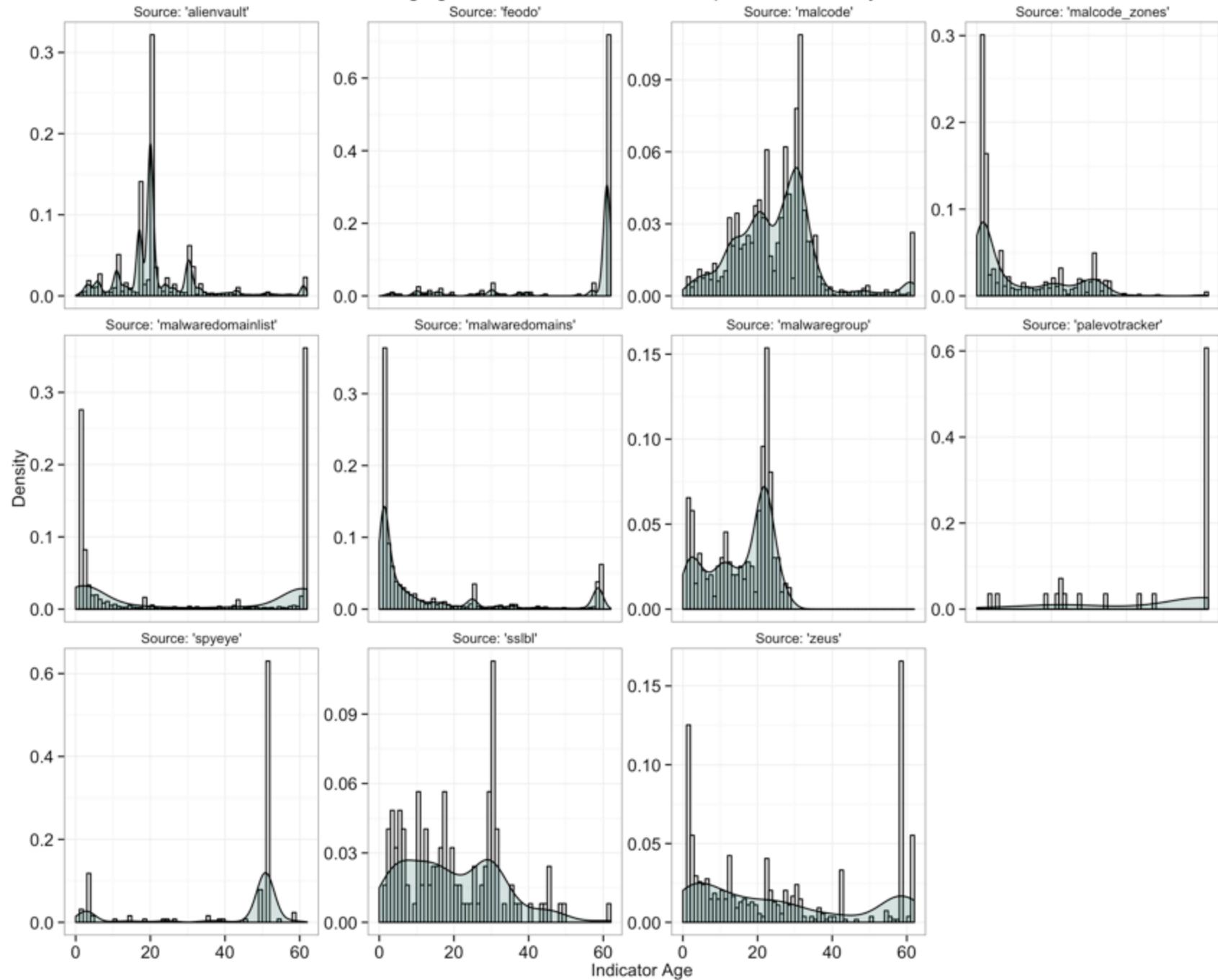
```
##   country conf.int.start conf.int.end    p.value
## 1:     DE      -0.002366      0.003411  0.7553
```

Aging Test – Is someone cleaning  
this up eventually?

### Aging Test - Inbound Data - Sampled Time: 61 days



### Aging Test - Outbound Data - Sampled Time: 61 days



# Uniqueness Test



# Uniqueness Test

- “Domain-based indicators are unique to one list between 96.16% and 97.37%”
- “IP-based indicators are unique to one list between 82.46% and 95.24% of the time”

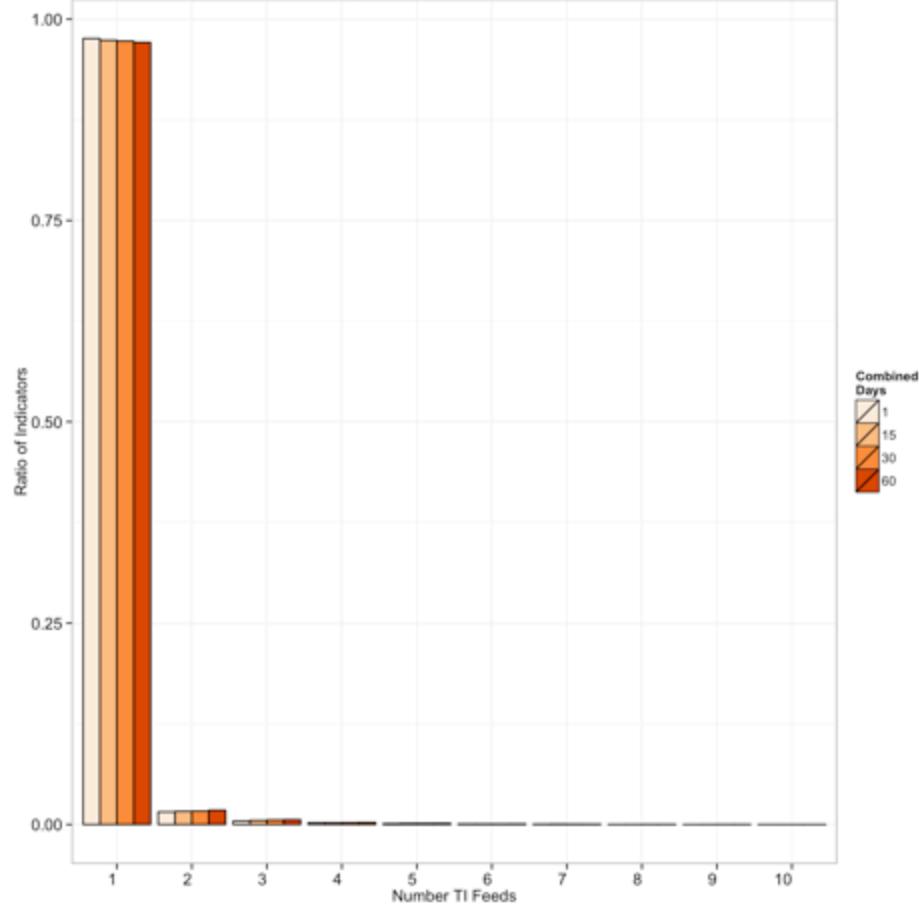


## Blacklist Ecosystem Analysis Update: 2014

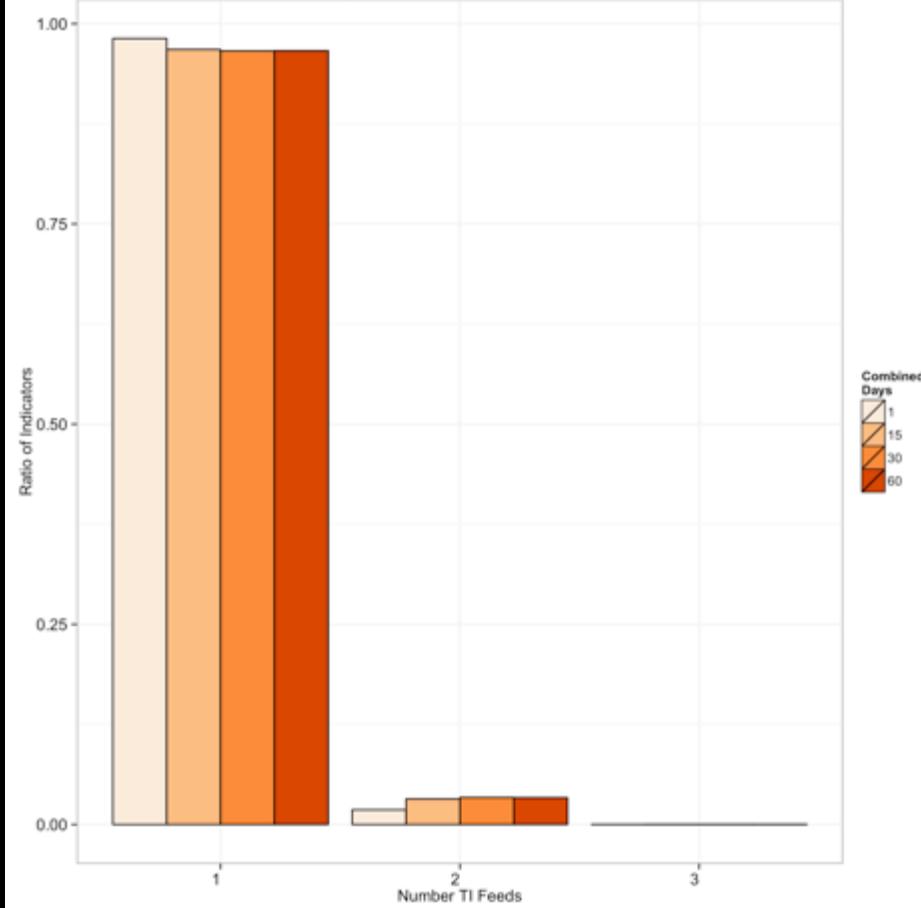
Leigh Metcalf, Jonathan M. Spring  
CERT® Division, Software Engineering Institute  
Carnegie Mellon University  
[netsa-contact@cert.org](mailto:netsa-contact@cert.org)  
Publication CERTCC-2014-82

December 2014

Uniqueness Test - Inbound Data



Uniqueness Test - Outbound Data



```
##      count      ratio days
## 1:      1 0.9759170     1
## 2:      1 0.9737684    15
## 3:      1 0.9727630    30
## 4:      1 0.9710713    60
```

```
##      count      ratio days
## 1:      1 0.9815816     1
## 2:      1 0.9678728    15
## 3:      1 0.9660537    30
## 4:      1 0.9663163    60
```

try some delicious and healthy snacks  
at our concession stand

# Intermission



# OPTION 1: Cool Story, Bro!

- Some of you are probably like:
  - “You Data Scientists and your algorithms, how quaint.”
  - “Why aren’t you doing some useful research like nation-state attribution?”



# OPTION 2: How can I use this awesomeness on my data?



# Use Case: Comparing Private Feeds

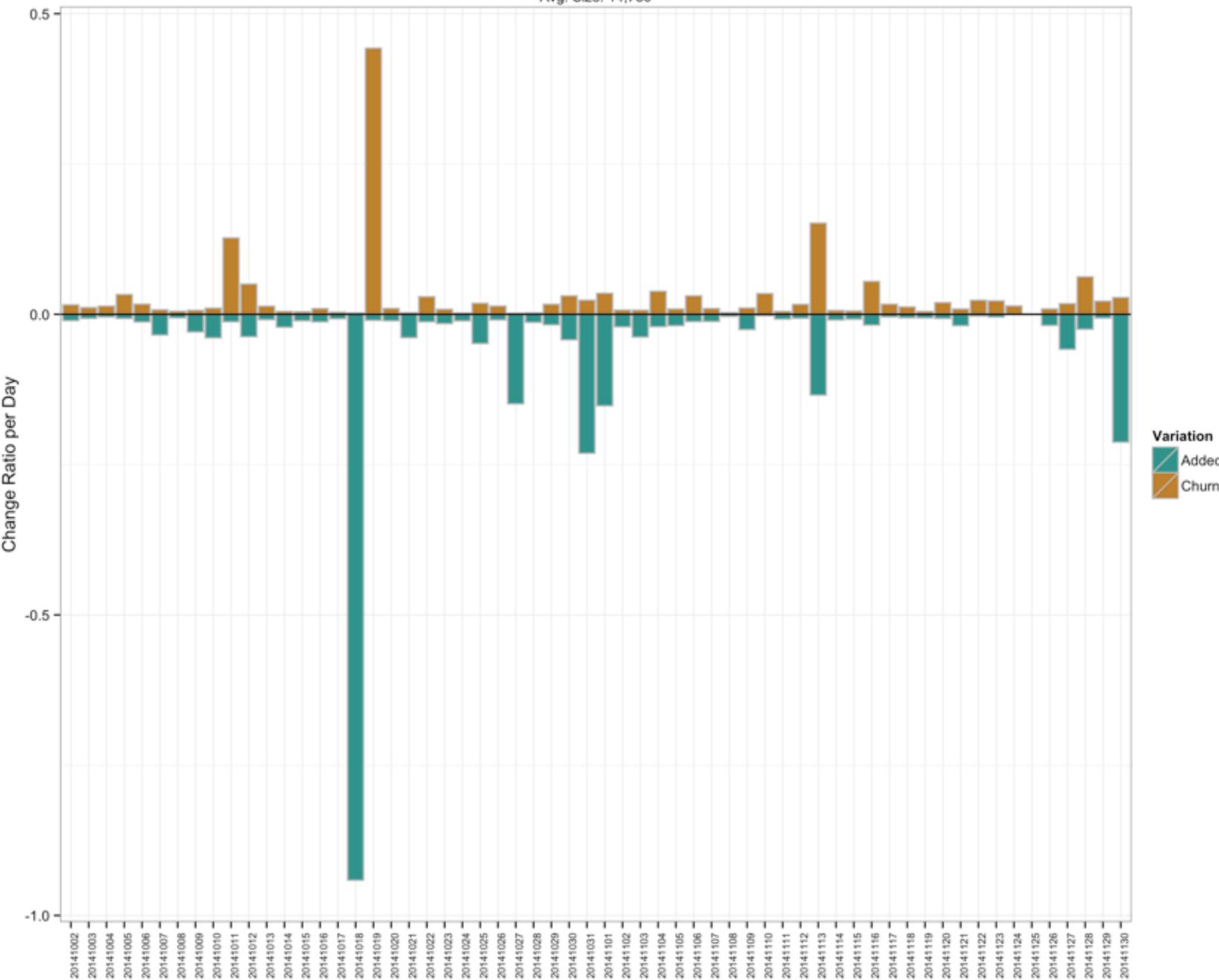
- How about using TIQ-TEST to evaluate a private intel feed?
- Trying stuff before you buy is usually a good idea. Just sayin'
- Let's compare a new feed, "private1", against our combined outbound indicators



## TIQ Novelty Test

Source Name: public\_outbound

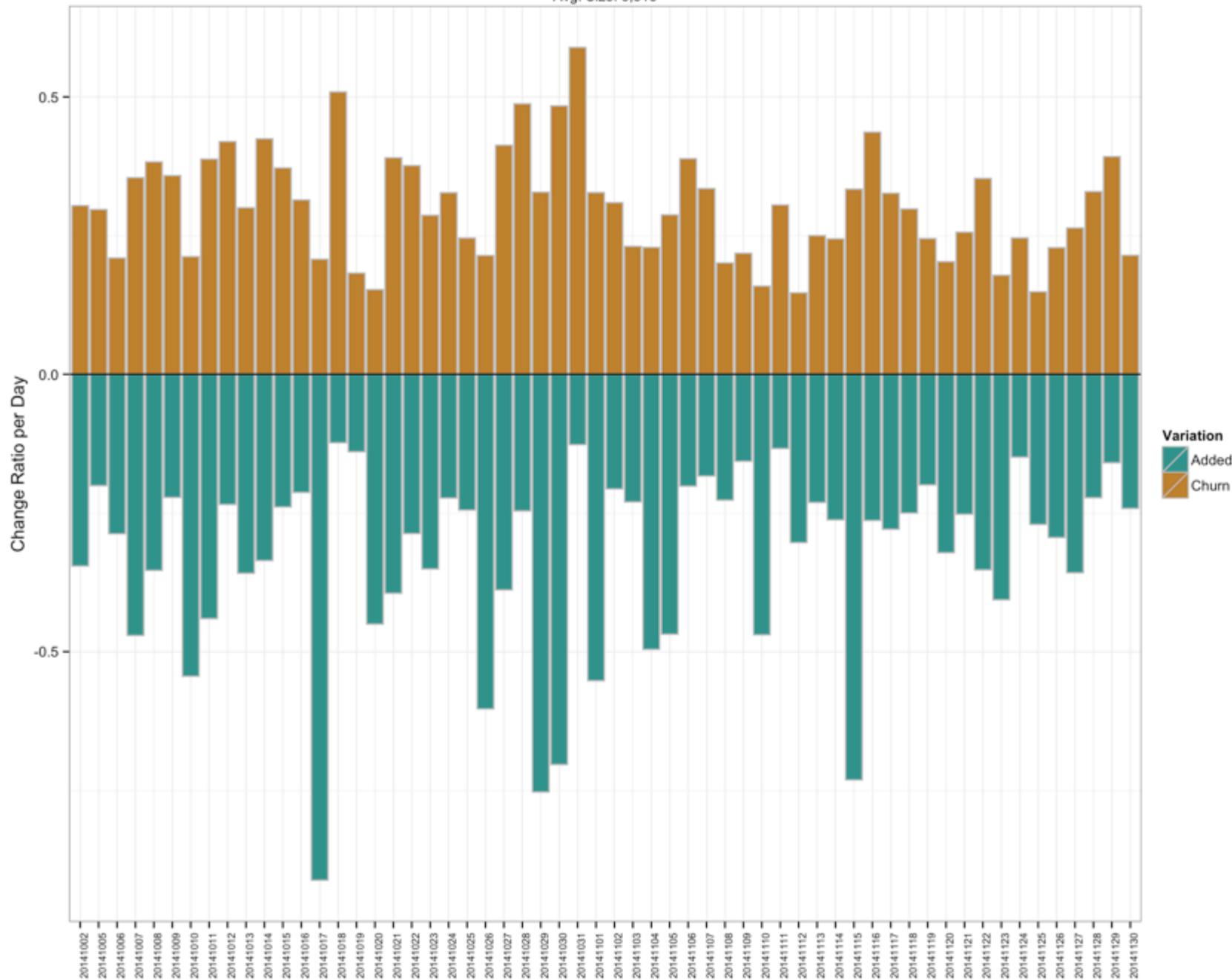
Avg. Size: 14,786



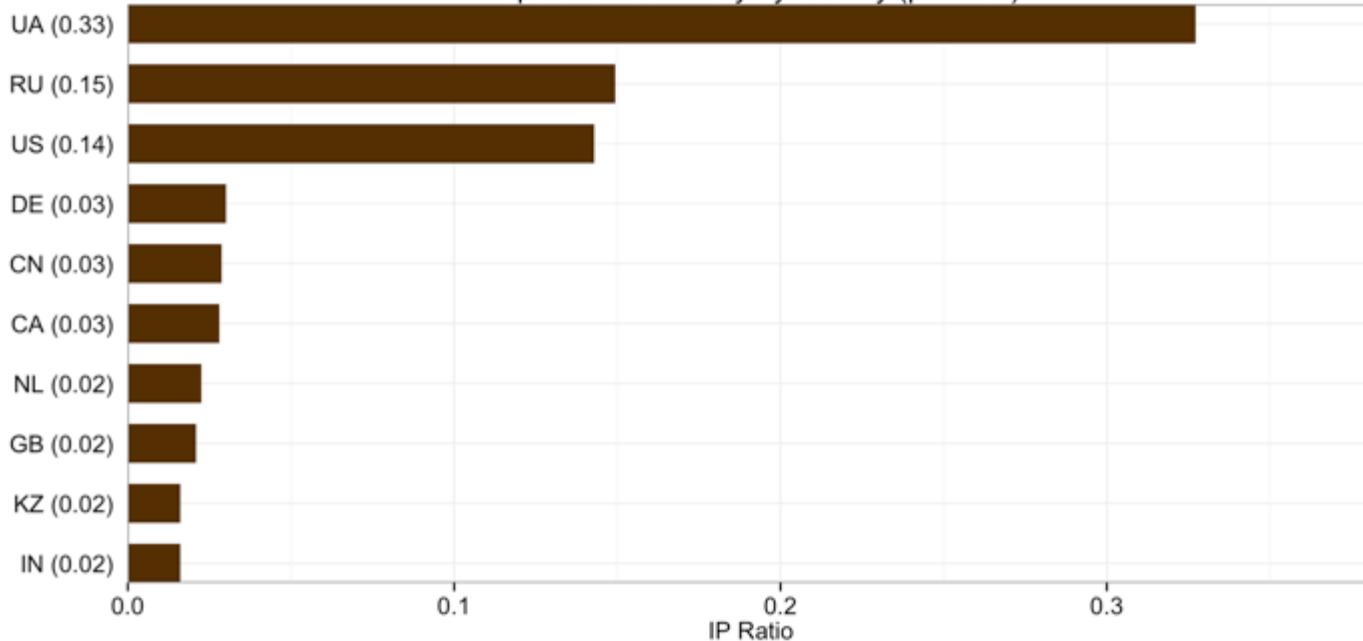
# TIQ Novelty Test

Source Name: private1

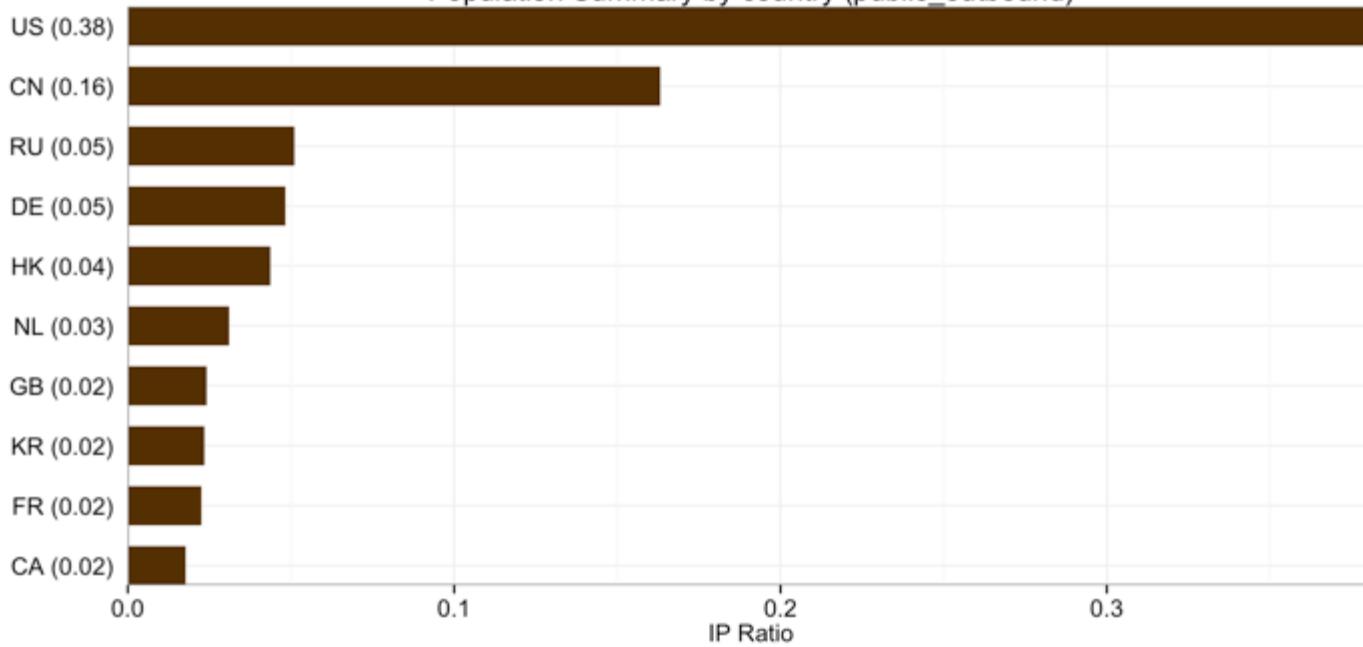
Avg. Size: 3,518



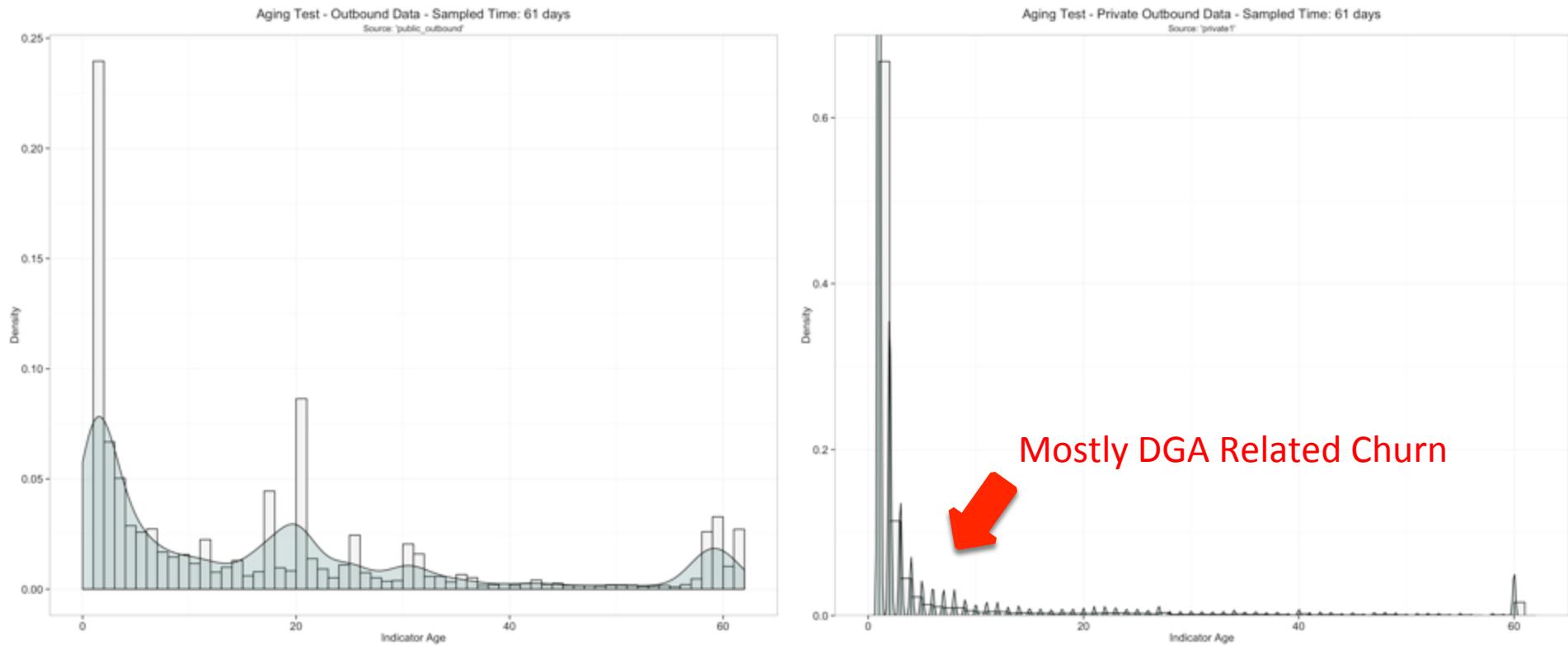
Population Summary by country (private1)



Population Summary by country (public\_outbound)

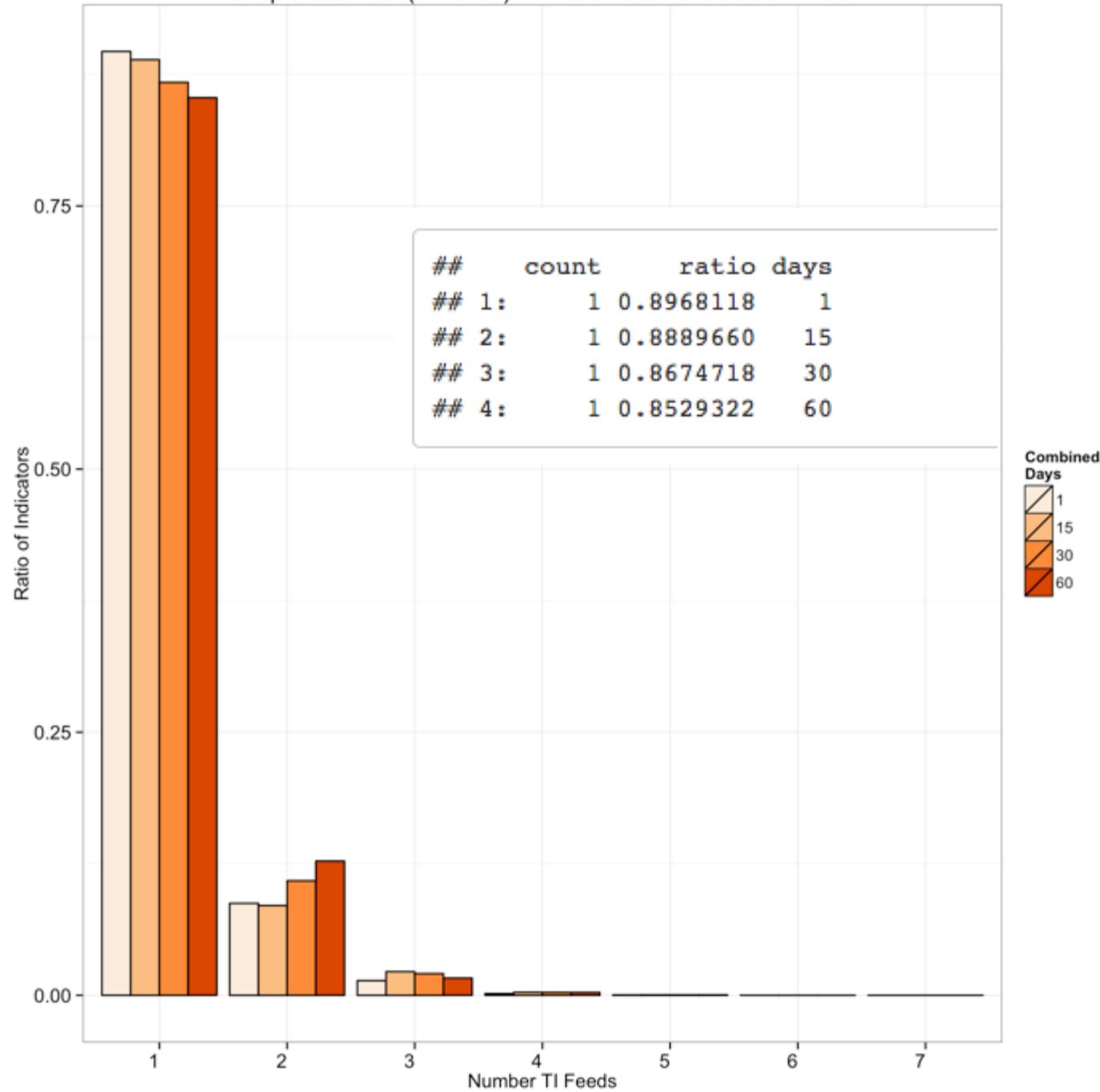


# Aging Test

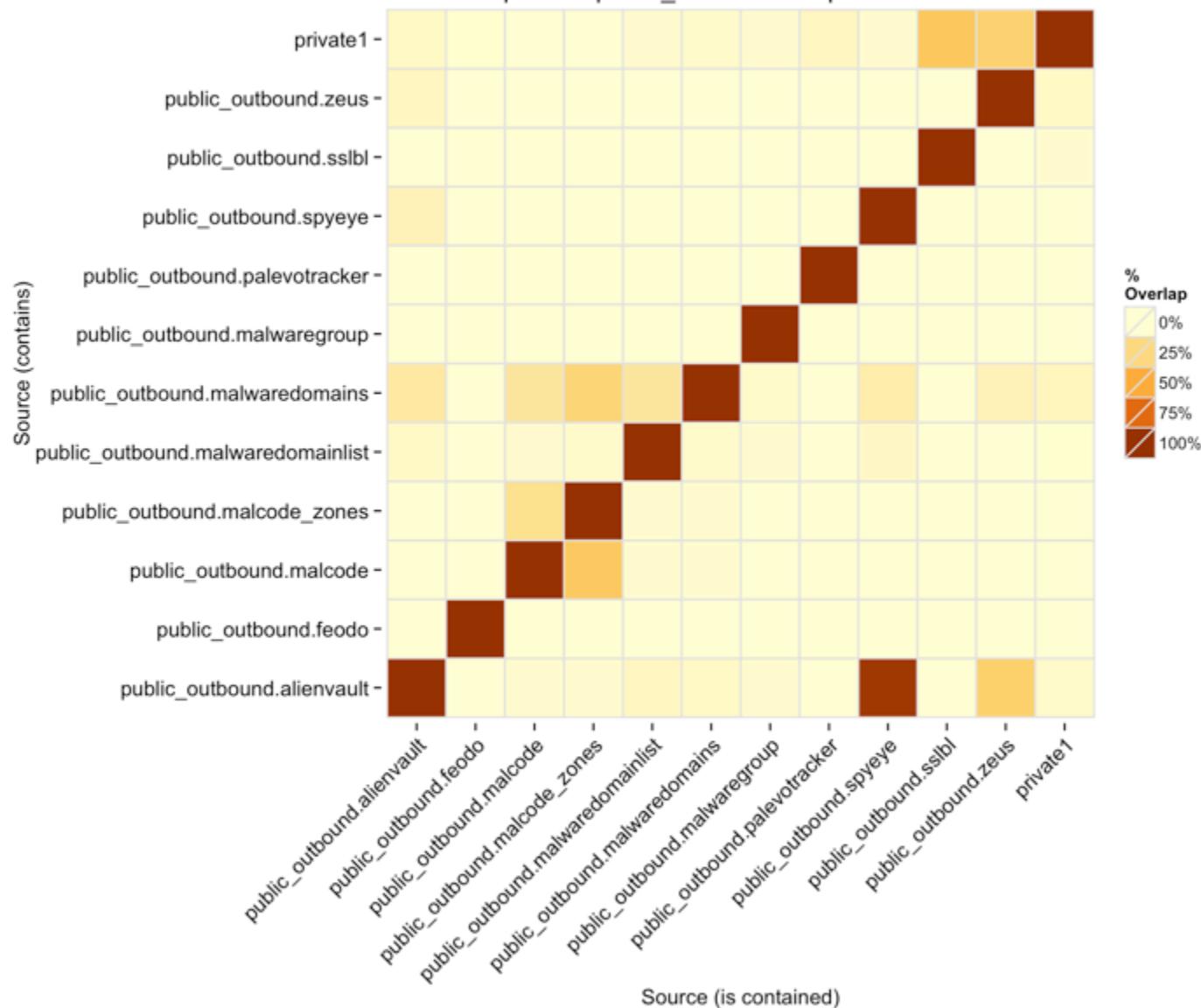


- I guess most DGAs rotate every 24 hours, right?
- Rotation means the private data is still “fresh”, from research or DGA generation procedures

Uniqueness Test (enriched) - Private Data vs. Outbound Data



Overlap Test - public\_outbound VS private1 - 20141101



1.

## A++

1. Relatively poor eBay feedback, compared to the community-accepted standard of **A+++++**. Often used for negative feedback responses or as revenge.
2. A++ is also a programming language rumored to be used by more than three people worldwide.

*"After getting my address from eBay, this seller broke into my house and killed my pets. A++"*

*"Man, this buyer sucks! I'll get him back by leaving him a mere A++ feedback."*

by **Jin64** September 02, 2007

A++ WOULD  
THREAT INTEL  
AGAIN

# Take Away

- Analyze your data. Extract value from it!
- Try before you buy! Different test results mean different things to different orgs.
- Try the sample data, replicate the experiments:
  - <https://github.com/mlsecproject/tiq-test-Winter2015>
  - <http://rpubs.com/alexcpsec/tiq-test-Winter2015>

# Greets and Thanks!

- @kafeine and John Bambenek for access to their feeds
- @kylemaxwell, @paul4pc and all contributors for COMBINE
- @bfist for his work on <http://sony.attributed.to>
- @hrbrmstr for IPEW and contribs for TIQ-TEST
- All the MLSec Project community peeps!



**Your gift of a few contributions  
Can help a starving data  
scientist.**

Healthcare (62)

Information  
(51)

Finance (52)

Educational  
(61)

Public (92)

Retail  
(44-45)

- Q&A?
- Feedback!

Alex Pinto  
@alexcpsc  
@MLSecProject  
@NiddelCorp

Alexandre Sieira  
@AlexandreSieira  
@MLSecProject  
@NiddelCorp



"The measure of intelligence is the ability to change."  
- Albert Einstein