



Google Analytics and How to Avoid Bad Data

An enterprise study of data quality



Brian Clifton - Author | co-founder Verified-Data.com

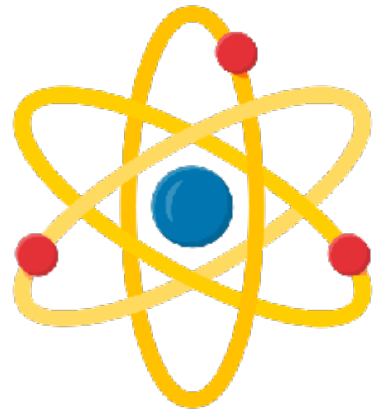




About Brian Clifton



Circa 1974...!



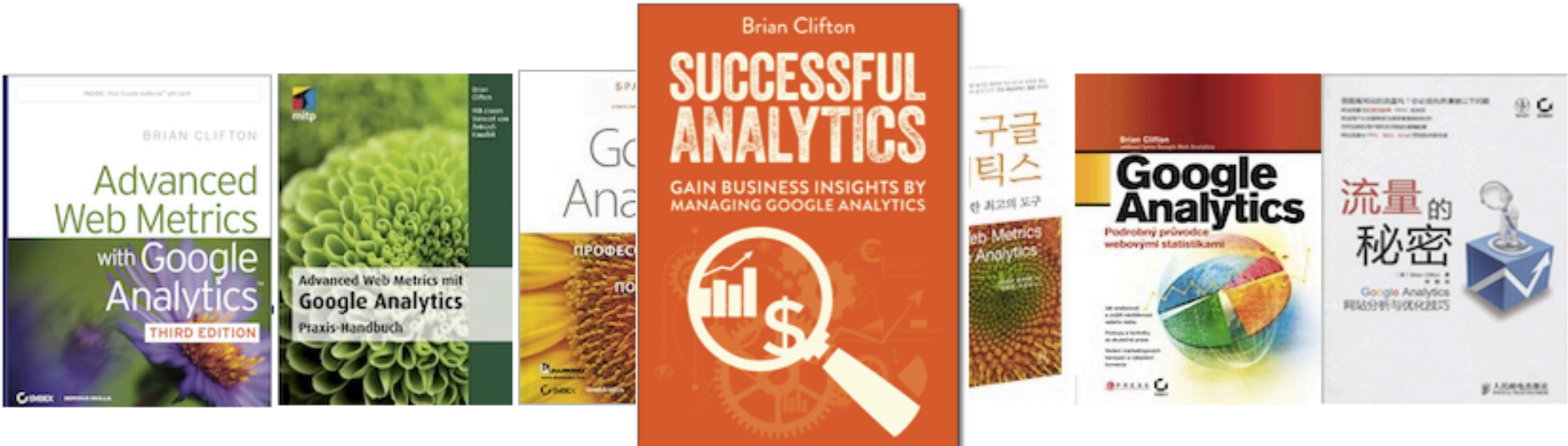
Trained as a scientist
(Chemistry)



Head of Web Analytics
EMEA 2005-8



Automated audit tool



100,000 copies sold

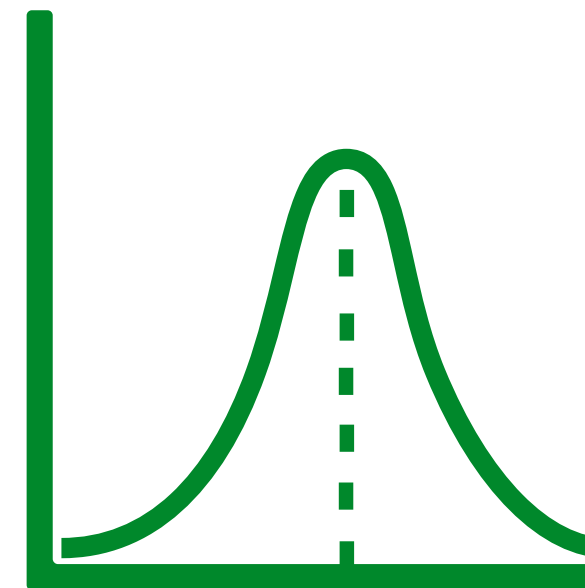
A black and white photograph of a rock climber standing on a high, rocky ledge. The climber is wearing a helmet, a light-colored long-sleeved shirt, and dark pants. They are holding a rope that extends down to another person who is hanging from the rope. The background shows a vast, rugged landscape with a deep canyon and a forested valley below. The sky is clear and blue.

Do you TRUST
your data?

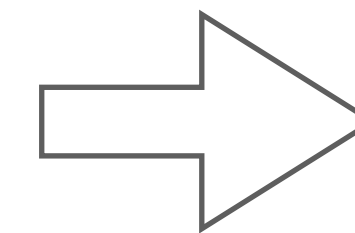
Why this lack of trust...?

- **Setup** is often not understood
- Collection is **rarely verified** - often "*smells*"
- **Poor governance**
 - particularly Google Analytics

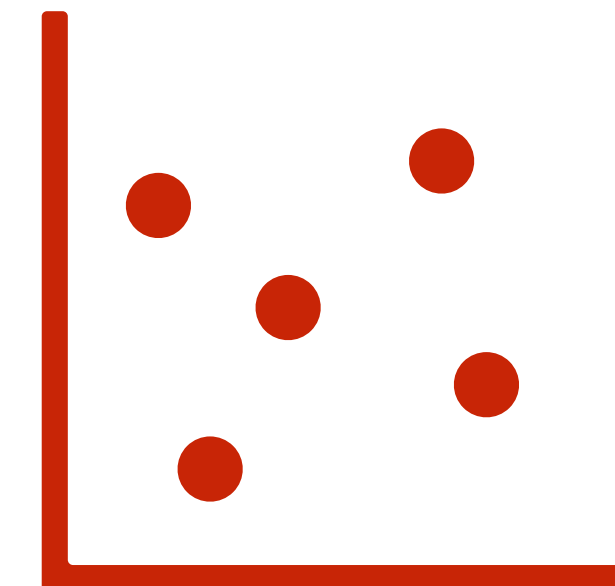
What you want



Clean data

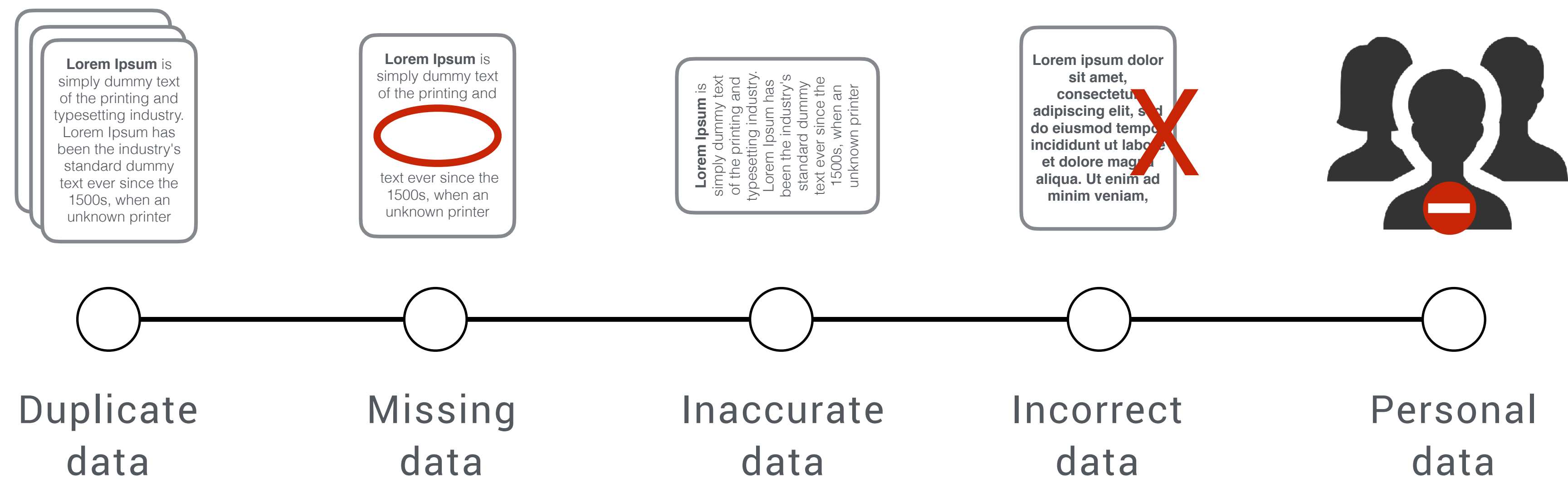


What you end up with



Poor quality

BAD DATA IS...



But its **hard** to find because...

- Websites constantly change
- Changes usually performed by "others"
- Time pressures - campaign "laziness"
- **Bad data looks just like good data!**

"Needles in impossible haystacks."



Good Data?

Bad Data?

The Data Quality Study

(the method)

Weight scorecard technique

Governance

Accuracy

Scorecard

	Weight	Website 1	Website 2	Website 3	Website 4
Core Principals					
Account Structure	1.0	!	!	!	!
Data Validity	1.0	X	X	!	!
Error Page Tracking	1.0	X	✓	✓	✓
Deployment & Coverage	1.0	!	!	✓	X
Privacy & Compliance					
PII Compliance	2.0	X	✓	X	X
Privacy & Consent Compliance	1.0	✓	✓	✓	✓
Cookie Report	1.0	!	!	✓	!
Marketing & Engagement					
Google Ads Tracking	1.0	X	X	X	X
Advanced Visitor Segmentation	1.0	✓	X	✓	✓
Campaign Tracking	1.0	X	X	!	✓
Site Search Tracking	1.0	X	!	!	!
File Download Tracking	1.0	X	✓	!	X
Outbound Link Tracking	1.0	!	✓	✓	!
Event Tracking	1.0	!	X	!	!
Goal Setup	1.0	✓	X	X	X
E-commerce					
Transaction Tracking	2.0	X	X	!	✓
Quality Index		29.4	40.7	53.6	39.3

Summary of over 200 data quality tests

Scorecard

	Weight	Website 1	Website 2	Website 3	Website 4
Core Principals					
Account Structure	1.0	!	!	!	!
Data Validity	1.0	X	X	!	!
Error Page Tracking	1.0	X	✓	✓	✓
Deployment & Coverage	1.0	!	!	✓	X
Privacy & Compliance					
PII Compliance	2.0	X	✓	X	X
Privacy & Consent Compliance	1.0	✓	✓	✓	✓
Cookie Report	1.0	!	!	✓	!
Marketing & Engagement					
Google Ads Tracking	1.0	X	X	X	X
Advanced Visitor Segmentation	1.0	✓	X	✓	✓
Campaign Tracking	1.0	X	X	!	✓
Site Search Tracking	1.0	X	!	!	!
File Download Tracking	1.0	X	✓	!	X
Outbound Link Tracking	1.0	!	✓	✓	!
Event Tracking	1.0	!	X	!	!
Goal Setup	1.0	✓	X	X	X
E-commerce					
Transaction Tracking	2.0	X	X		
Sum and normalise		29.4	46.4	46.4	46.4

Section score

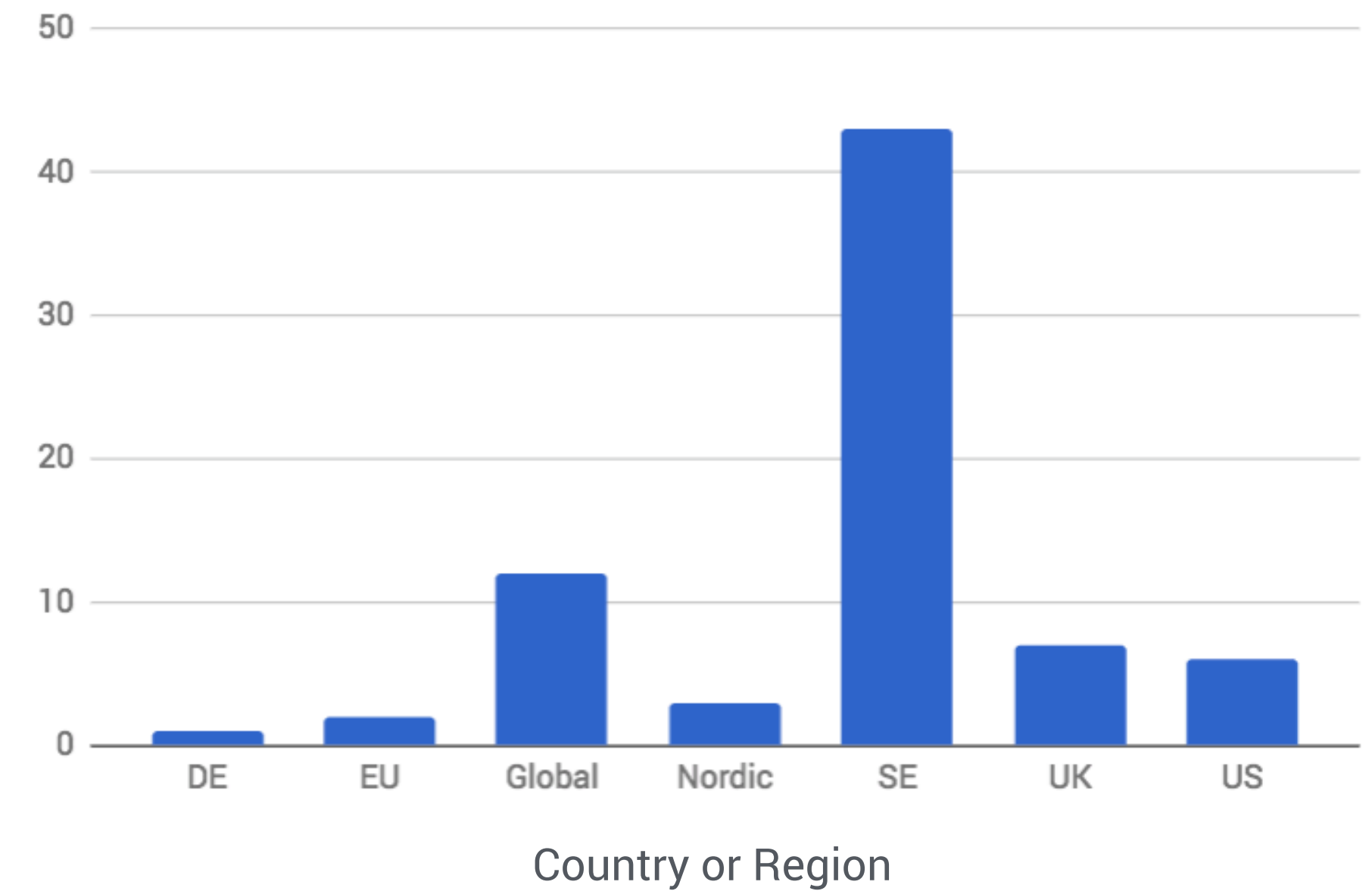
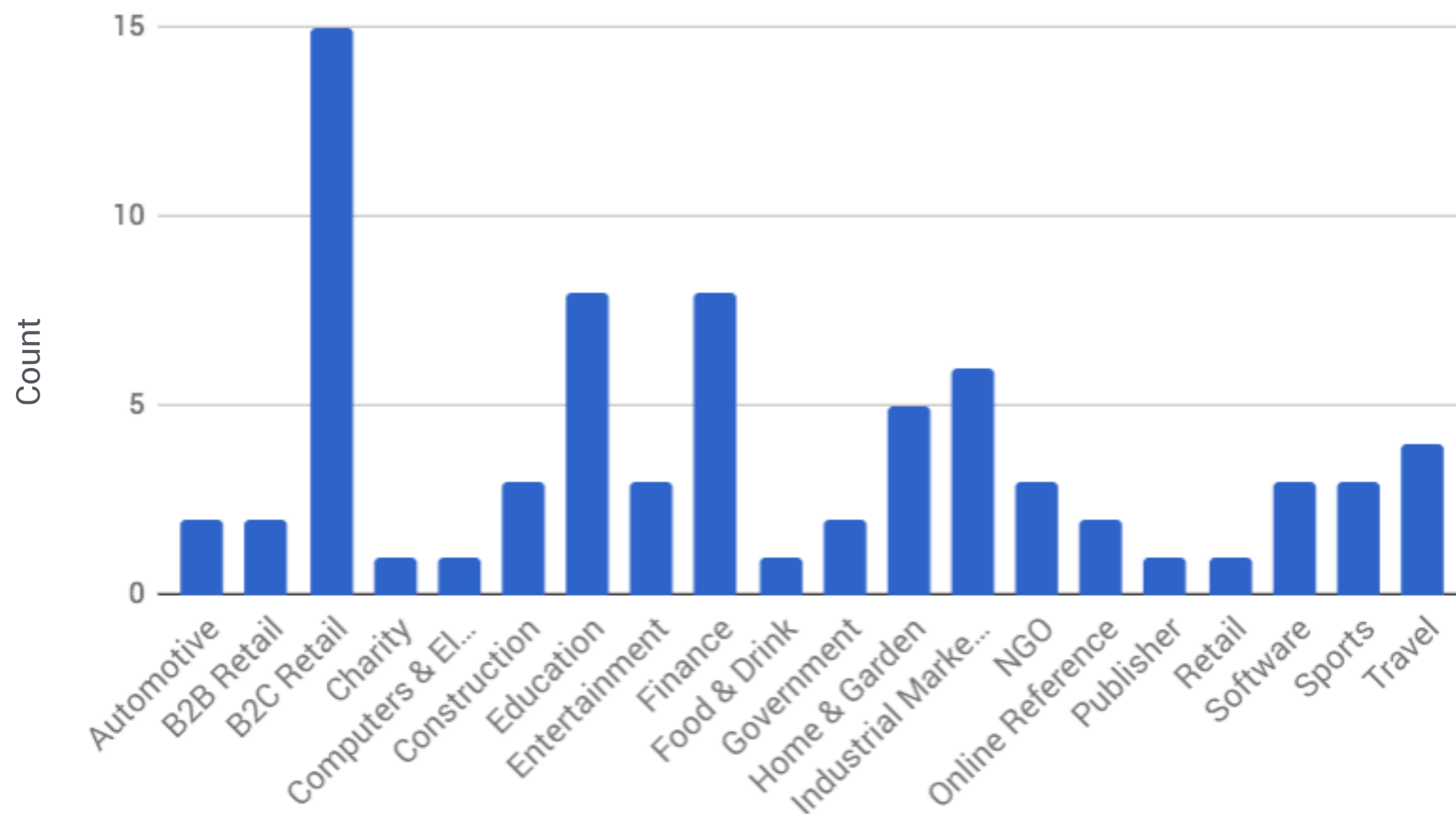
✓ = 10

! = 5

X = 0

Quality Index score out of 100

- All **brand leaders**, 13 with a global presence.
- **41%** had an e-commerce facility.
- Monthly visits range from 100k - 100 million (three > 100M).



The Results

(a visualising data quality)

Scorecard

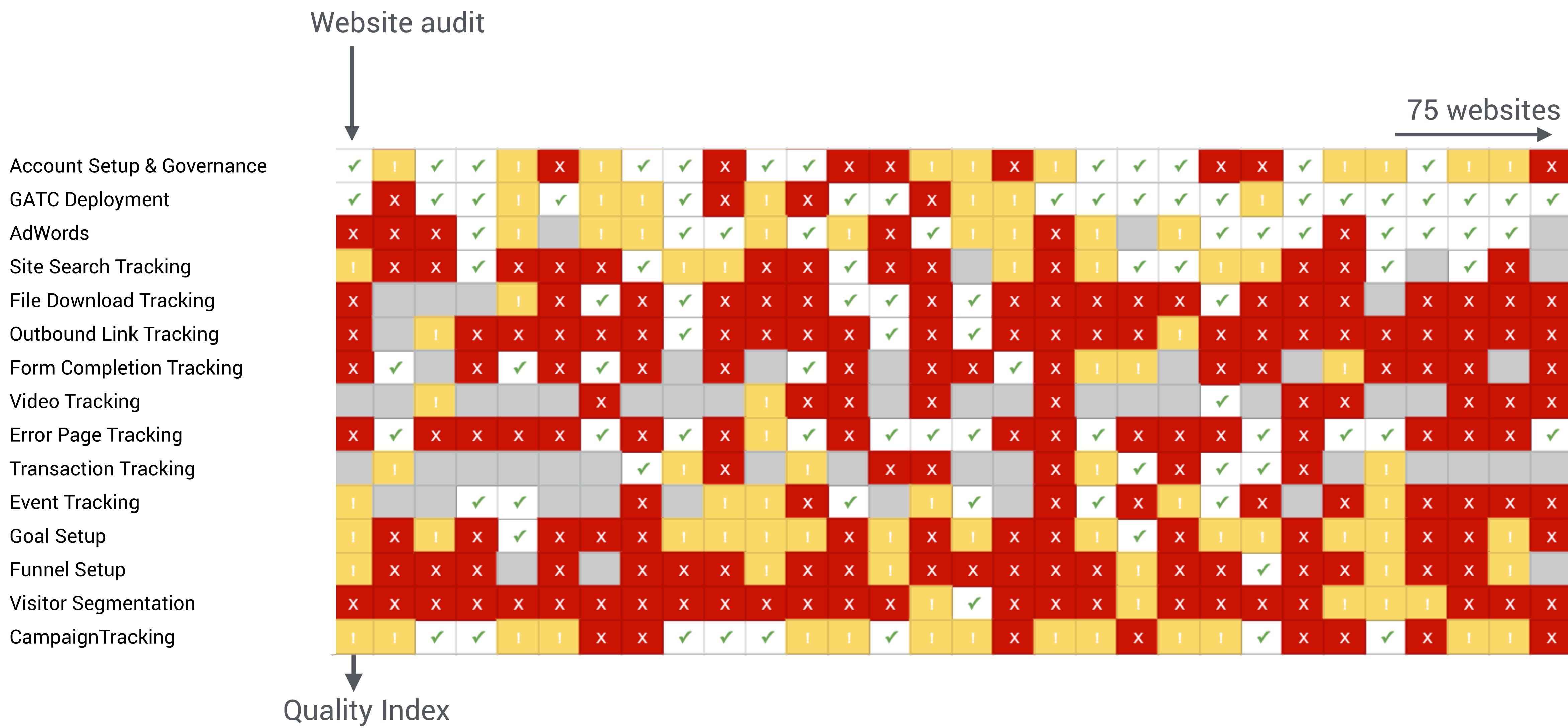
	Weight	Website 1	Website 2	Website 3	Website 4
Core Principals					
Account Structure	1.0	!	!	!	!
Data Validity	1.0	X	X	!	!
Error Page Tracking	1.0	X	✓	✓	✓
Deployment & Coverage	1.0	!	!	✓	X
Privacy & Compliance					
PII Compliance	2.0	X	✓	X	X
Privacy & Consent Compliance	1.0	✓	✓	✓	✓
Cookie Report	1.0	!	!	✓	!
Marketing & Engagement					
Google Ads Tracking	1.0	X	X	X	X
Advanced Visitor Segmentation	1.0	✓	X	✓	✓
Campaign Tracking	1.0	X	X	!	✓
Site Search Tracking	1.0	X	!	!	!
File Download Tracking	1.0	X	✓	!	X
Outbound Link Tracking	1.0	!	✓	✓	!
Event Tracking	1.0	!	X	!	!
Goal Setup	1.0	✓	X	X	X
E-commerce					
Transaction Tracking	2.0	X	X	!	✓
Quality Index		29.4	40.7	53.6	39.3

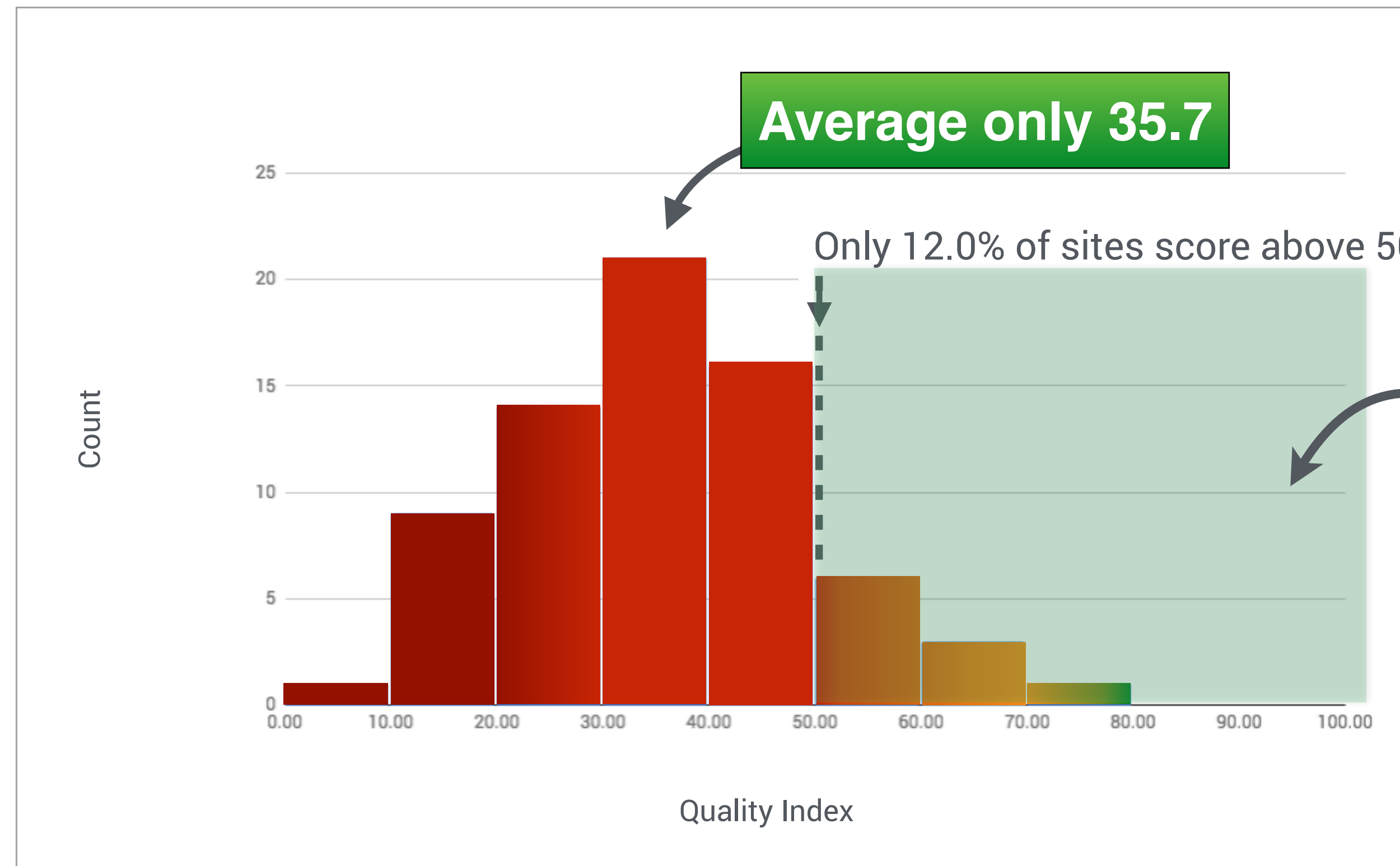
75 enterprise websites





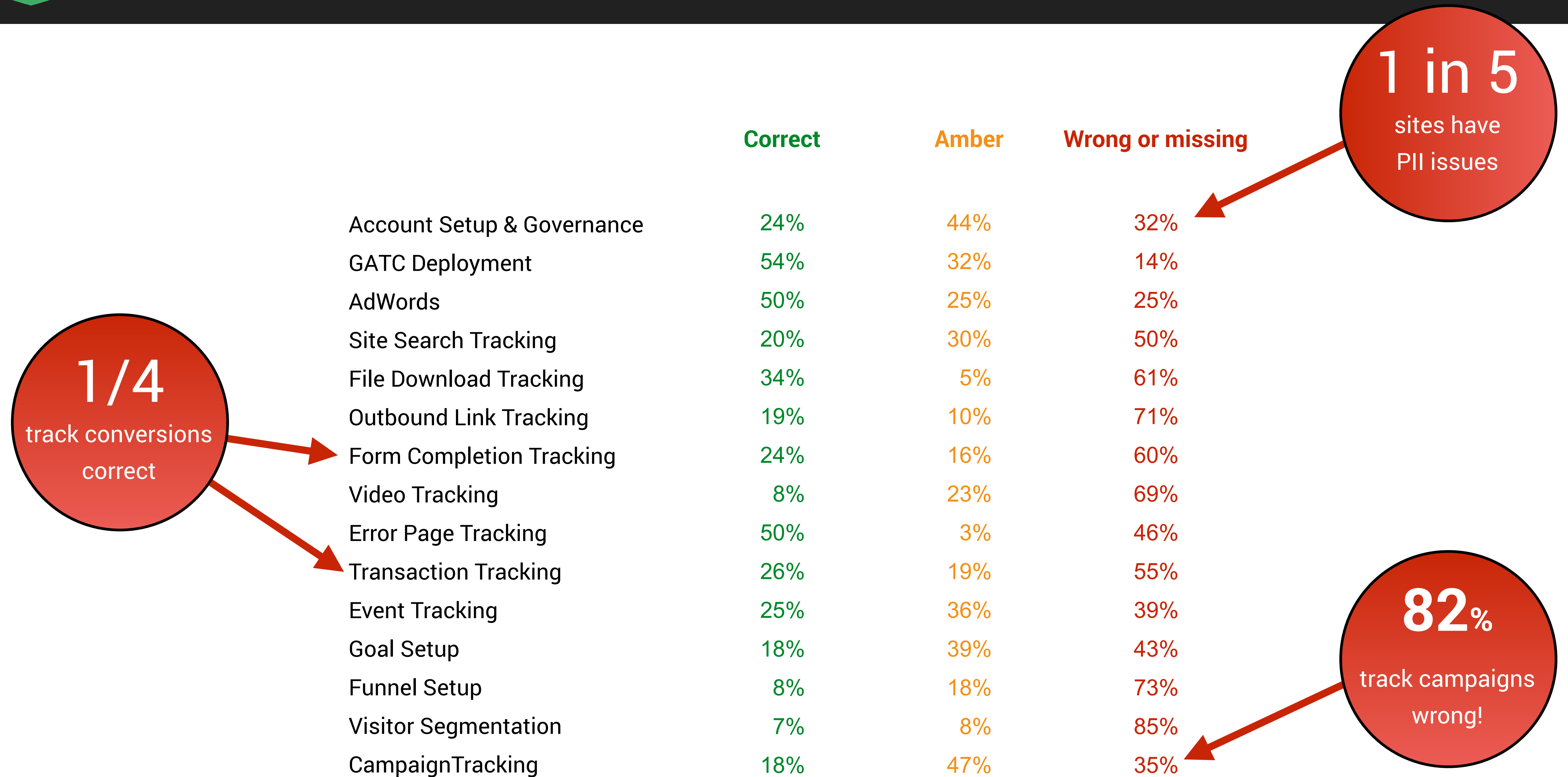
Data quality patterns - massive failures...!





Lowest = 4.5

Highest = 73.1



Examples of data quality issues

1 in 5
sites have
PII issues

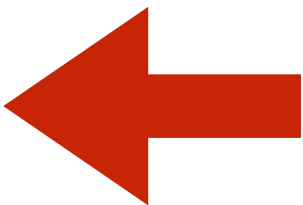




PII most common in **URLs** and **page titles**...

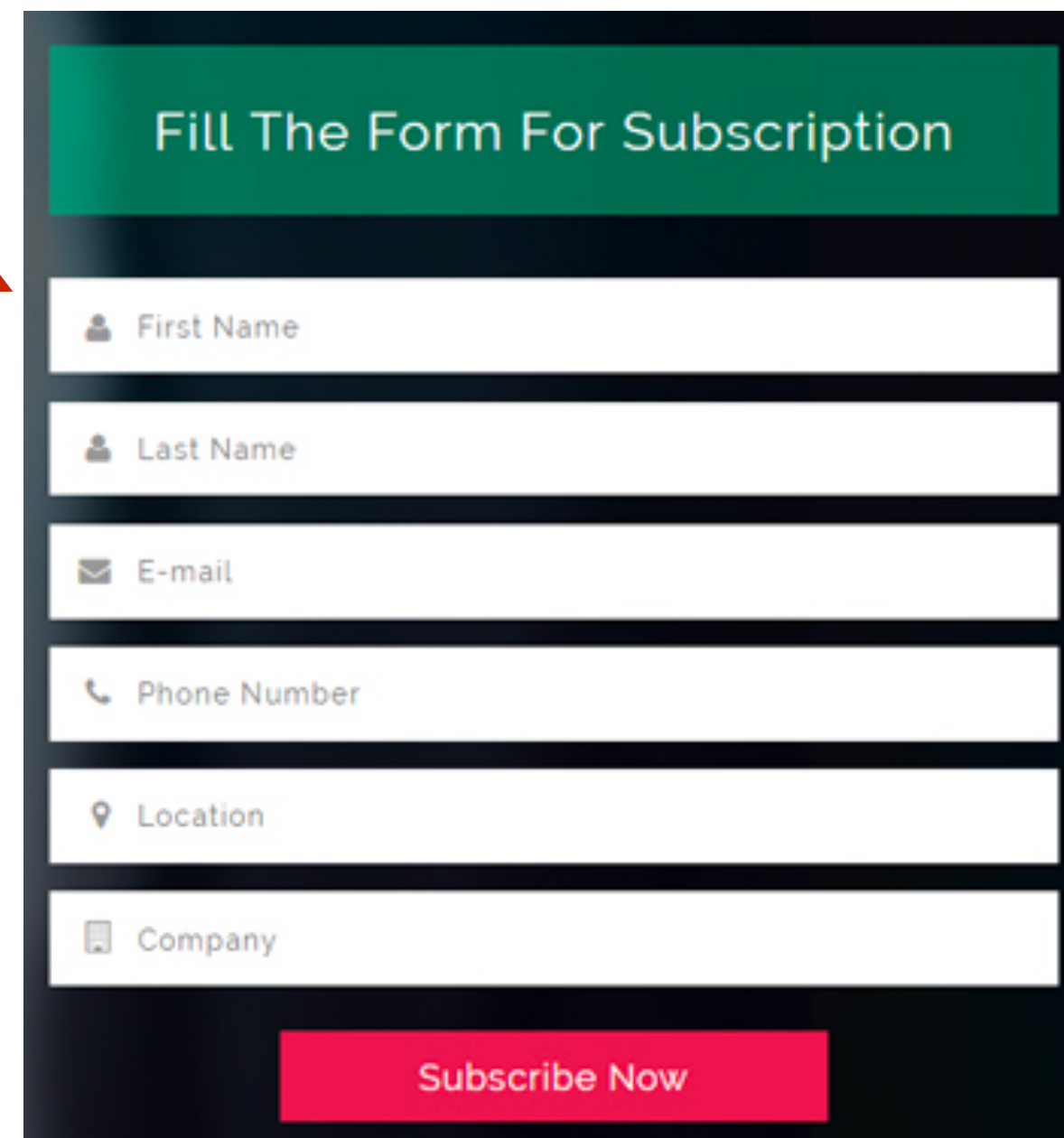
Usually by accident...

Page ?	Page Views ?
	13,382 % of Total: 0.10% (13,216,720)
1. /existing-appointment/log-in/?ins=wt16&email=ben_l i-8890@yahoo.com&password=daxiMgjxG4	59 (0.44%)
2. /existing-appointment/log-in/?ins=wt16&email=d aol.com&password=ngJ3UlifE7	45 (0.34%)
3. /existing-appointment/log-in/?ins=wt16&email=robbr @hotmail.com&password=zMVMbBheN7	44 (0.33%)
4. /existing-appointment/log-in/?ins=wt16&email=ched hotmail.co.uk&password=QF5HoNwzK8	40 (0.30%)
5. /existing-appointment/log-in/?ins=wt16&email=a ng@gmail.com&password=DxCe4plcl5	36 (0.27%)
6. /existing-appointment/log-in/?ins=wt16&email=mcco hotmail.co.uk&password=RLGt8tXnA9	33 (0.25%)
7. /existing-appointment/log-in/?ins=wt16&email=shini hotmail.com&password=9vVnbQ7oZ8	33 (0.25%)
8. /existing-appointment/log-in/?ins=wt16&email=TON YSECURITYGROUP.CO.UK&password=et7eNSyl F4	33 (0.25%)



Where does "accidental" PII come from...?

Via forms submitted using GET method...



Fill The Form For Subscription

First Name

Last Name

E-mail

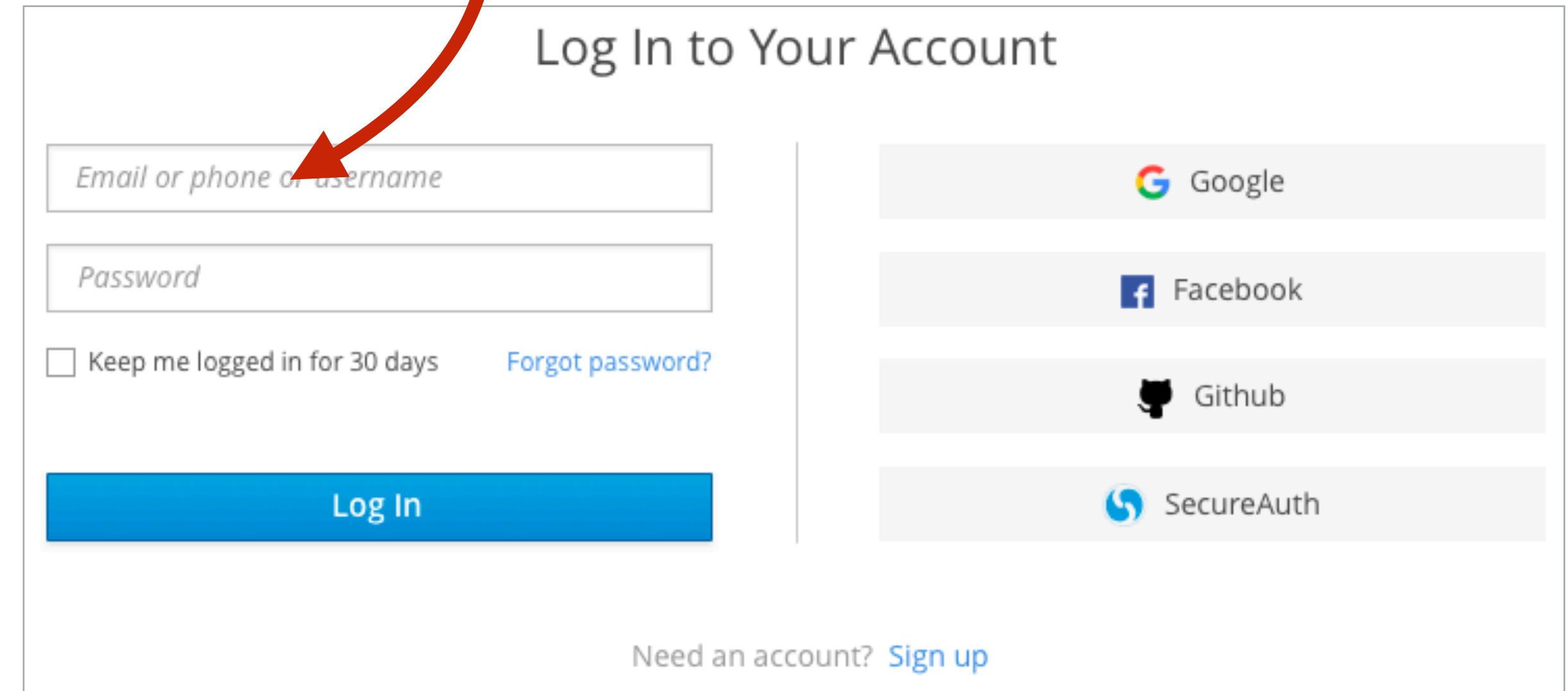
Phone Number

Location

Company

Subscribe Now

<title> tags: "Welcome Brian@..."



Log In to Your Account

Email or phone or username

Password

☐ Keep me logged in for 30 days [Forgot password?](#)

Log In

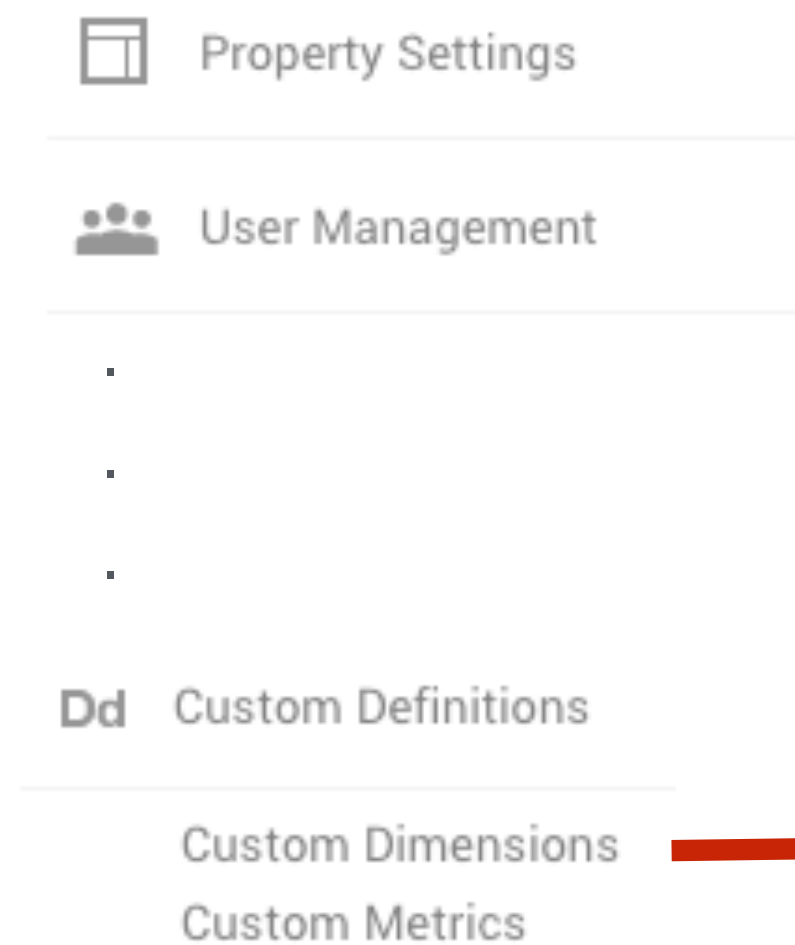
Google

Facebook

Github

SecureAuth

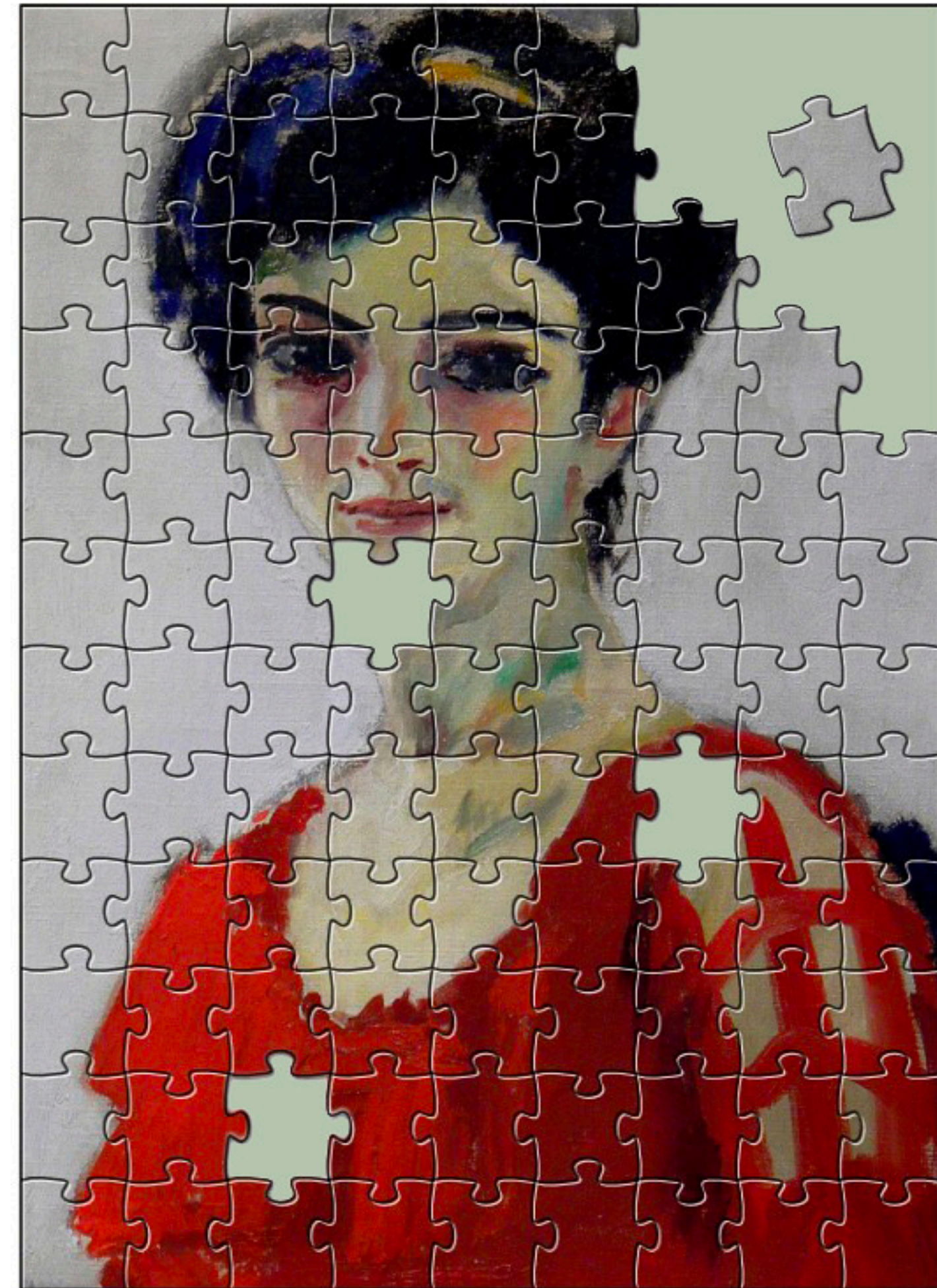
Need an account? [Sign up](#)



Per field, this is not PII... BUT

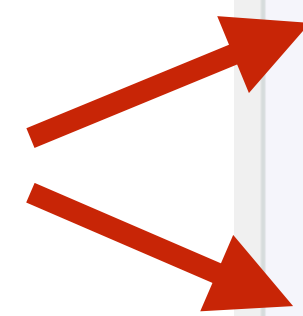
UI gender
UI postal code
UI city
UI state
UI birthday

Jigsaw effect of data triangulation



Only deletion is available
for Google Analytics and
its very blunt...

ALL URLs removed!



Create a new data deletion request

Property ID
198011909

Data deletion requests are run in UTC. [Learn more](#)

Start date *
10/1/2019

End date *
10/31/2019

Fields to delete *
URL

[Submit](#)

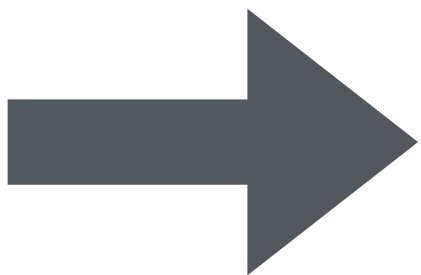
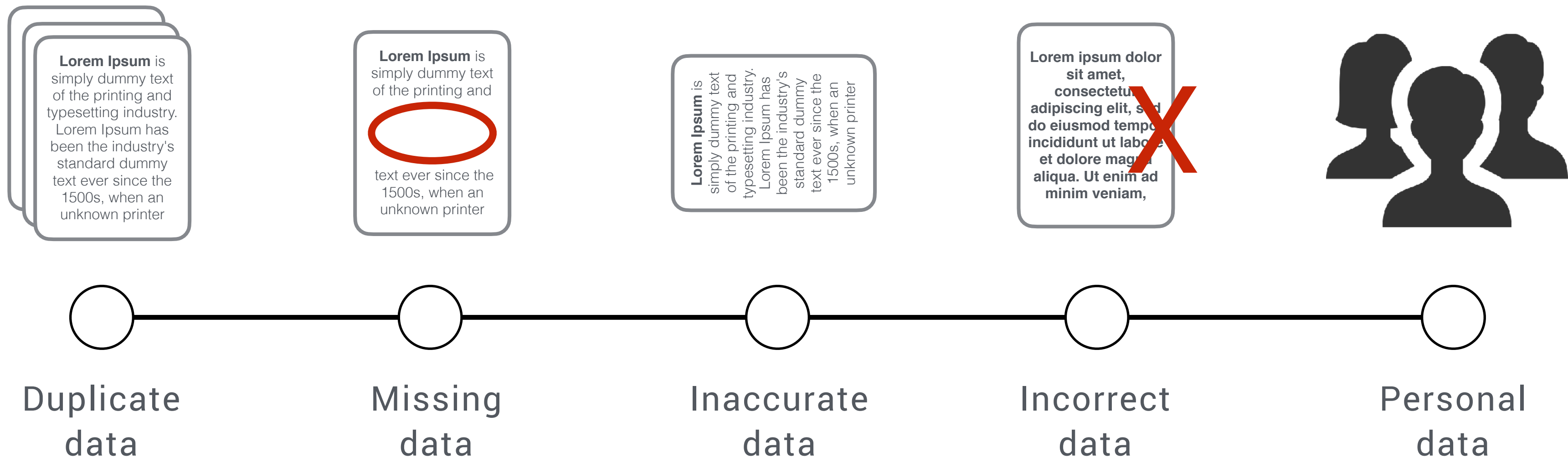
Only

1/4

Track transactions
correct!

Tran\$action Tracking

BAD DATA IS...

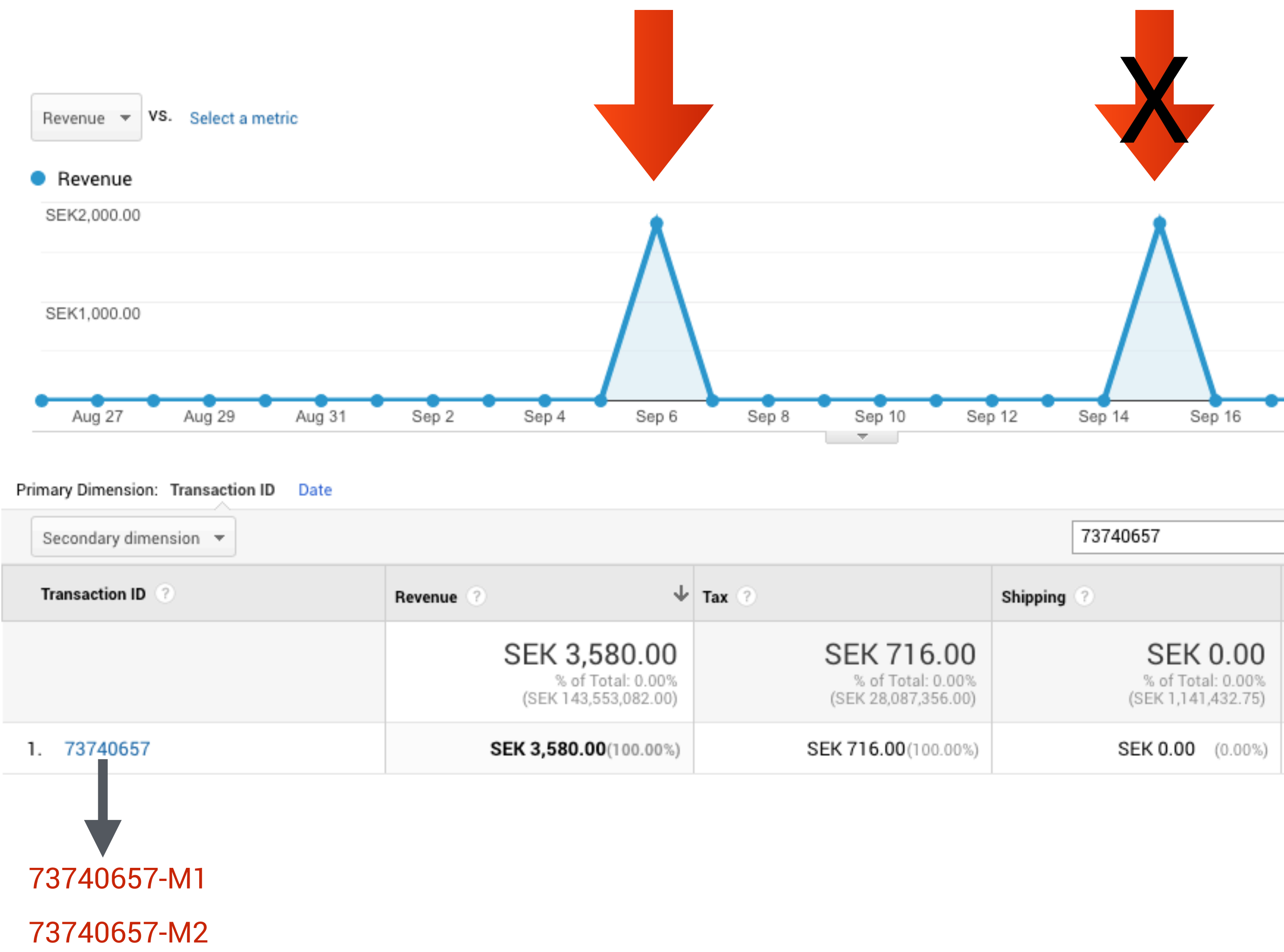


- Duplicate transactions
- Product lists, Promotions
- Inconsistent dataLayer
- Timing issues
- Name, address etc. (in the wrong place!)

The most common problem - Duplicate Transactions

If orders are not unique
your attribution will be a
mess....!

Overwrites the original referrer



If existing orders do need
updating, **append** to the
transID

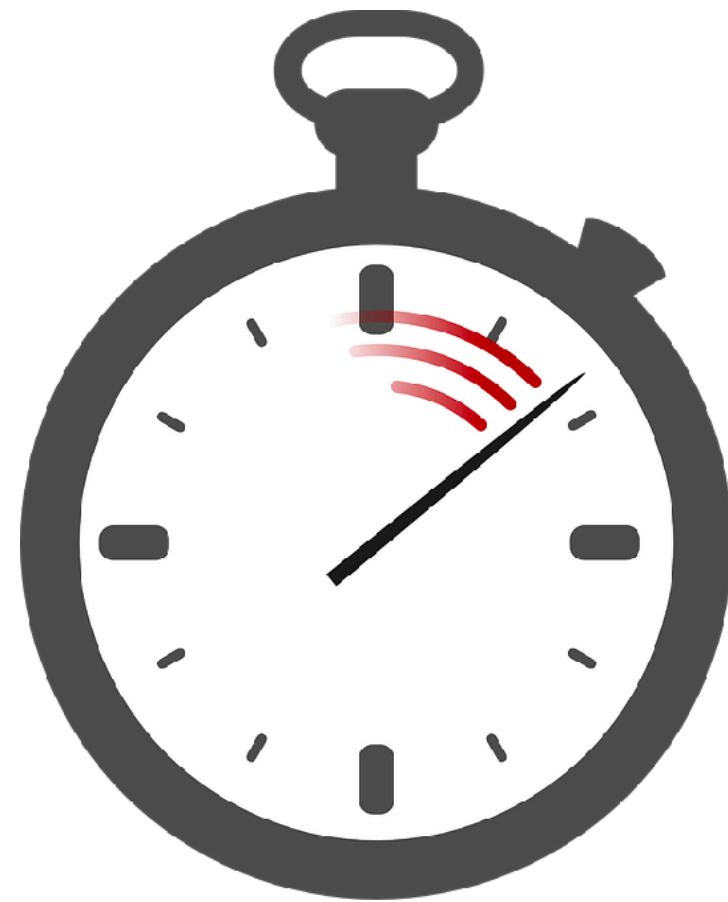
???

Transaction ID ?	Revenue ?	Tax ? ↑
	£1,024,352.00 % of Total: 0.71% (£143,553,082.00)	-£504,054.00 % of Total: -1.79% (£28,087,356.00)
500. 73740036	£98.00 (0.01%)	-£78.00 (0.02%)
501. 73739525	£118.00 (0.01%)	-£76.00 (0.02%)
502. 73764528	£99.00 (0.01%)	-£76.00 (0.02%)
503. 73703352	£5,890.00 (0.57%)	-£73.00 (0.01%)
504. 73788966	£1,890.00 (0.18%)	-£72.00 (0.01%)
505. 73742391	£89.00 (0.01%)	-£71.00 (0.01%)
506. 73700533	£399.00 (0.04%)	-£70.00 (0.01%)
507. 73720309	£648.00 (0.06%)	-£70.00 (0.01%)
508. 73788793	£149.00 (0.01%)	-£70.00 (0.01%)
509. 73780313	£669.00 (0.07%)	-£66.00 (0.01%)

Personal v B2B sales!
Did not match backend process.

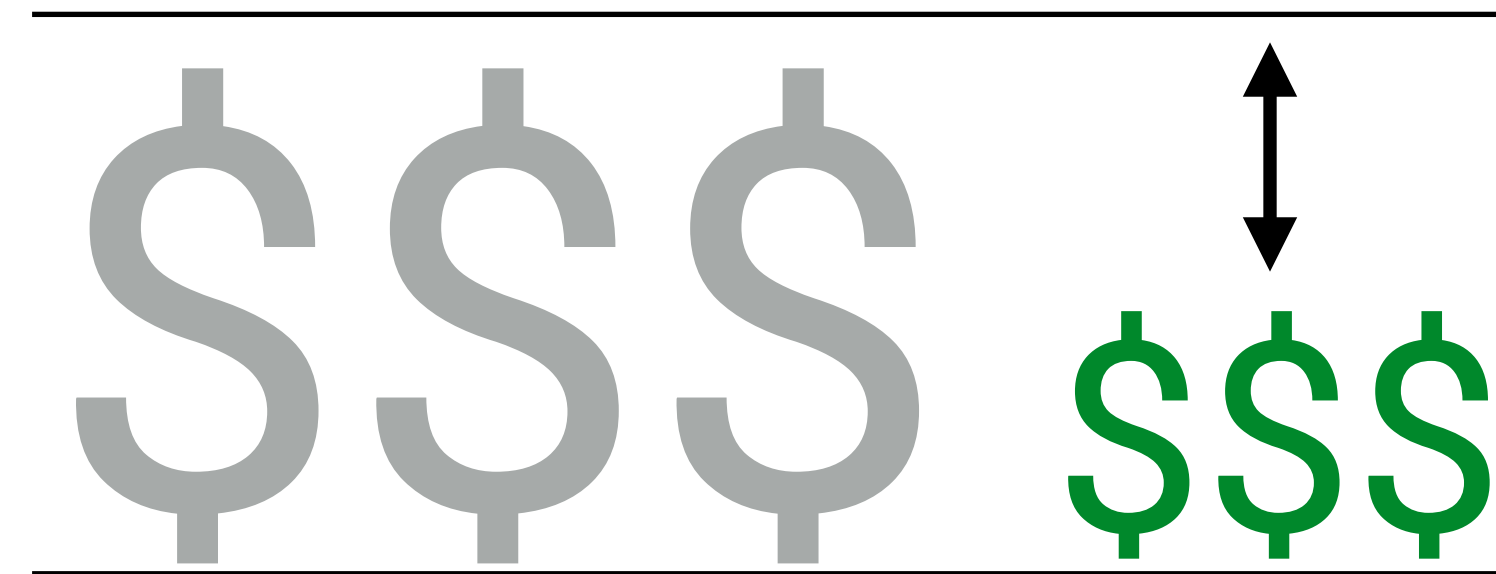
Lessons learned:
KIS; Don't make me think

TIMING



- Google Analytics captures in **real time**.
- Backend processing is usually **batched**.

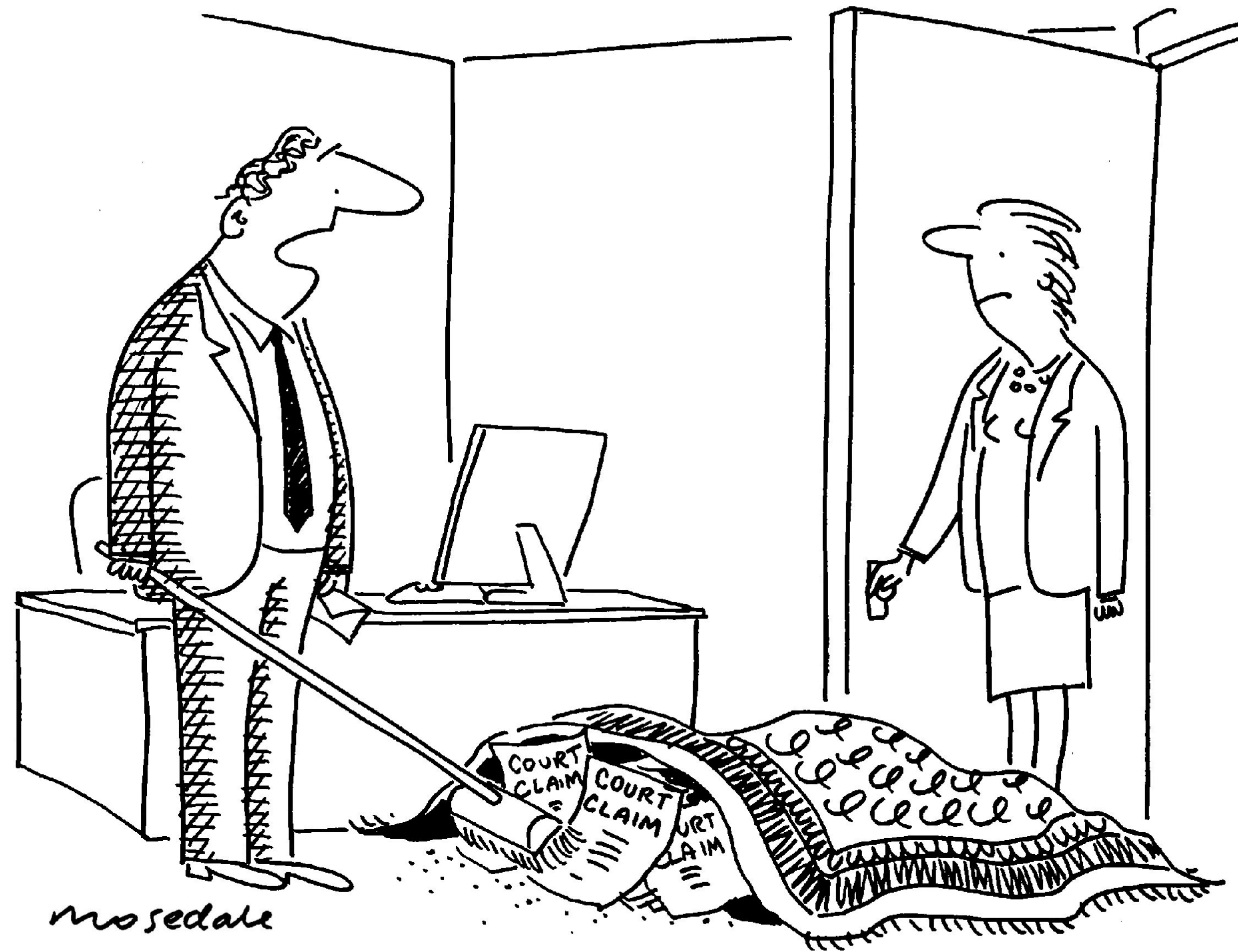
REFUNDS



- **Avoid refunds** in web analytics.
- Note them for backend reconciliation (or use a separate GA property).

**Why are we in this
position...?**

Hard to find issues are **easy to not know about** (or even ignore)...



After all, these are needles in impossible haystacks...

This mess can be avoided...



- **Stop** auditing manually
 - Or you will lose good staff!
- **Automate** the heavy lifting
 - Retain/recruit good staff
- **Fix** the priorities
 - Get your Quality Index score 50+
- **Monitor** regularly
 - Keep your QI > 80

If you have \$100 to “*make smart decisions using data*”,
invest **\$1** to monitor and verify its quality.