# Protecting Data Sources, Protecting Personal Data
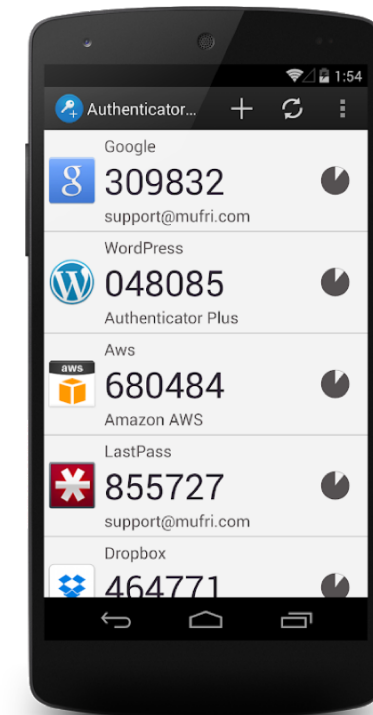
Simson L. Garfinkel
Senior Scientist, Confidentiality and Data Access
U.S. Census Bureau

4th Kavli Symposium
Tuesday, February 20
Austin, TX

# Raise your hand if you use two-factor authentication to protect your email account

# Two factor authentication

# Protecting Data Sources, Protecting Personal Data



Sources

Collection

Processing

Dissemination

*Communications Security*

*Storage Security*

*Publications Security*

# Outline

**Communications security**: How do you obtain confidential information from your sources?

**Storage security**: How do you maintain your secrets?

**Publication security**: How do you control the information released by your publication to prevent the inadvertent release of confidential information?

# Bio: Simson L. Garfinkel

1987  Freelance science writer    The Boston Globe
                                  WIRED

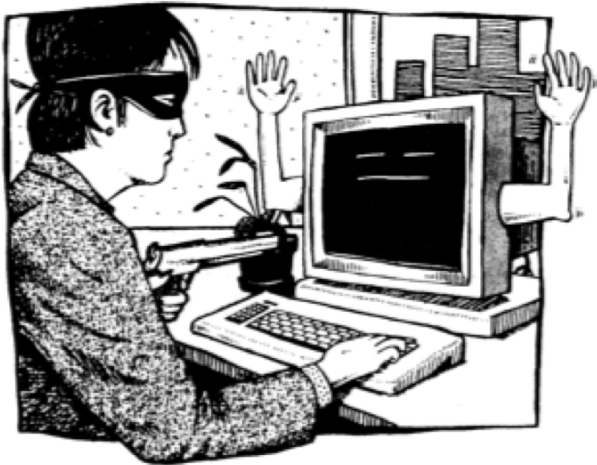2006  Associate Professor         NAVAL POSTGRADUATE SCHOOL

2015  Computer Scientist          NIST

2017  Computer Scientist          United States Census Bureau

# I have spent 29 years trying to secure computers...



**An Introduction to Computer Security**
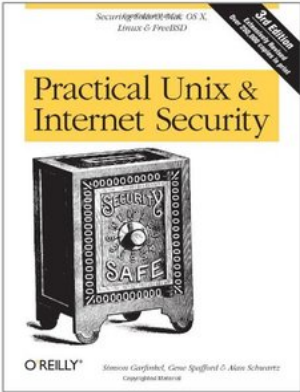[Part 1]

Simson L. Garfinkel

"Spies," "vandals," and "crackers" are out there, waiting to get into—or destroy—your databases.

L AWYERS MUST UNDERSTAND is-sues of computer security, both for the protection of their own inter-ests and the interests of their clients.

Lawyers today must automatically recognize insecure computer systems and lax operating procedures in the same way as lawyers now recognize
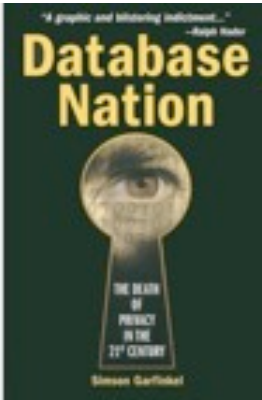
39

*The Practical Lawyer*
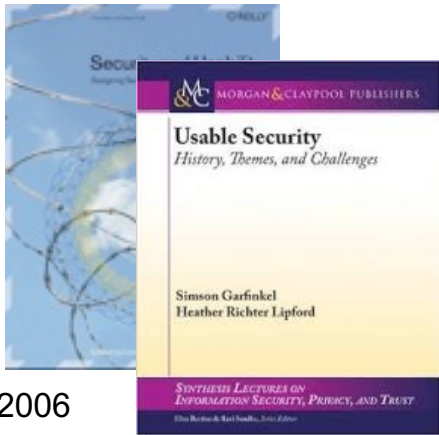Sept. 1987

System Security
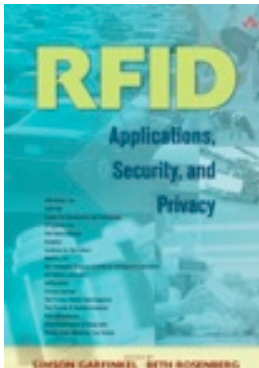


1991

Privacy Policy



2000

Usable Security



2006

2014

Internet of Things



2006

# Today's systems are frequently less secure than those of the 1970s.

**Poor security** is inherent in many information systems.

- Attack is easier and cheaper than defense.
- Cyber "defense in depth" does not work
  *a single vulnerability compromises.*
- It's easier to break things than to fix them.

**Network Connectivity** makes it easier to exploit vulnerable computers.
Fortunately, most journalists have modest security needs.

United States™ **Census** Bureau

U.S. Department of Commerce
Economics and Statistics Administration
U.S. CENSUS BUREAU
census.gov

# A methodology for thinking about your security needs

Identify your **critical assets** and **interactions** — what you are trying to protect.

Identify potential **threats** — what you are trying to protect against.

Identify potential **vulnerabilities** — how your threat could be harmed

Identify **risks** — the potential for harm

Asset + Threat + Vulnerability = Risk

# There are many risk equations

Asset + Threat + Vulnerability = Risk

Risk = Threat * Vulnerability * Consequences

Risk = Impact * Probability

These equations *shouldn't* be solved quantitatively.

United States Census Bureau

U.S. Department of Commerce
Economics and Statistics Administration
U.S. CENSUS BUREAU
census.gov

# Communications Security

# Communications with sources: Securing data in flight

## Primary risk: interception



Email



Phone



In-Person meeting

## Asset: content & reputation

# Which of these is has the most interception _risk_?



Answer depends on:

- **Threat** — who is attacking?

- **Vulnerability** — how are they attacking?

- **Consequence** — what is the impact of an interception?

# In-person meetings are risky

# The Acela Spy

**SLATE**

The shocking things I've learned by eavesdropping on Amtrak.

*By Amy Webb*

"On Amtrak, powerful people talk loudly and spill secrets.

"This is my conclusion based on five years' field research commuting on Amtrak's Acela between cities along the East Coast."

United States Census Bureau

U.S. Department of Commerce
Economics and Statistics Administration
U.S. CENSUS BUREAU
census.gov

Former NSA Director Michael Hayden was overheard on the phone to reporters on an Amtrak train. It happens more often than you might think.

# Eavesdropping email or phone requires *access.*

There are two points of access:

1. The end-point devices.

2. The network.



Primary threat:
spyware and malware



Primary threat:
interception

# Encryption doesn't protect against malware

"https:" encryption protects data in flight against interception.



**S/MIME and PGP** (message encryption) also protects data at rest. See NIST SP800-188, "Trustworthy Email."

# Storage Security

# Storage Security

## Local Storage



## Cloud Storage



Issues: Physical Access • Logical Access

# Most of the data crimes in recent years have been unauthorized access to stored data.

## Physical access:

- Attacker physically removed the data.

## Logical access:

- Computer system allowed access
- Data were not encrypted *to the attacker*.
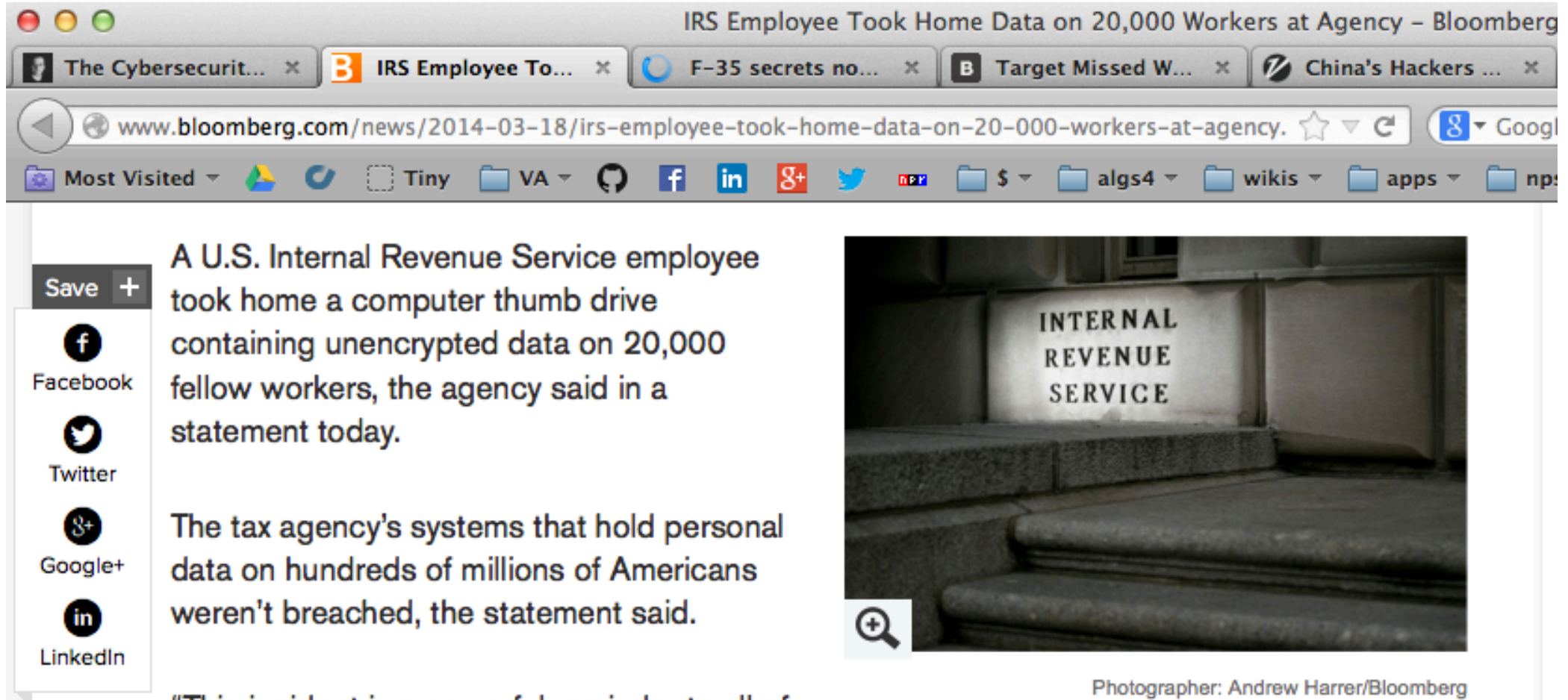
**Asset?**

**Vulnerability?**

**Threat?**

**Consequence?**

# May 2013: Edward Snowden steals millions of documents from the US National Security Agency

# March 2014:
# IRS Employee Took Home Data on 20,000 Workers



IRS Employee Took Home Data on 20,000 Workers at Agency – Bloomberg

www.bloomberg.com/news/2014-03-18/irs-employee-took-home-data-on-20-000-workers-at-agency.

A U.S. Internal Revenue Service employee took home a computer thumb drive containing unencrypted data on 20,000 fellow workers, the agency said in a statement today.

The tax agency's systems that hold personal data on hundreds of millions of Americans weren't breached, the statement said.

Photographer: Andrew Harrer/Bloomberg

# March 2014:
# Stolen F-35 secrets show up in China's stealth Fighter

# Sept. 2014: Celebrity photos stolen from iCloud

## Apple Admits Celebrity Photos Were Stolen In Targeted Hack

United States Census Bureau

U.S. Department of Commerce
Economics and Statistics Administration
U.S. CENSUS BUREAU
census.gov

# Protecting Local Storage

Physical security.

Disk encryption.

Off-site backups.



MacOS FileVault



Oakland CA fires, 1989

# Protecting Network Storage

Two-factor access

Account recovery

Google Drive          GMAIL          DropBox          OneDrive

# Example: Google Authenticator's 2-factor authentication protections against password stealing.

# Universal Second Factor (U2F)

United States Census Bureau

U.S. Department of Commerce
Economics and Statistics Administration
U.S. CENSUS BUREAU
census.gov

# Publication Security

# https://commons.wikimedia.org/wiki/ Commons:Photographs_of_identifiable_people



**Problems with attempts at anonymisation**

Black band over the eyes

Pixelated face

Cropped head

WORLD

# SECRETS OF HISTORY: The C.I.A. in Iran -- A special report.; How a Plot Convulsed Iran in '53 (and in '79)

By JAMES RISEN    APRIL 16, 2000

For nearly five decades, America's role in the military coup that ousted Iran's elected prime minister and returned the shah to power has been lost to history, the subject of fierce debate in Iran and stony silence in the United States. One by one, participants have retired or died without revealing key details, and the Central Intelligence Agency said a number of records of the operation -- its first successful overthrow of a foreign government -- had been destroyed.

# June 16, 2000: the *New York Times* publishes on its Web site a leaked secret CIA report on its website.

Report published as a PDF.

The Times had attempted to redact the names Iranians who had assisted.

The Times "redacted" by putting black boxes over the PDF.

Cryptome.org removed the black boxes and re-published.

http://cryptome.org/cia-iran.htm

Date: Mon, 19 Jun 2000 08:19:45 -0400

To: intelforum@his.com

From: John Young jya@pipeline.com

Subject: Re: Complete CIA history of 1953 Iranian coup posted by New York Times

The digital means the NY Times used to black out names of persons it was advised might be put at risk by publication failed to do the job properly. All the deletions are readable. The unredacted report shall be published shortly on cryptome.org.

The unexpected consequences of digital security are worth pondering.

Date: Tue, 20 Jun 2000 11:04:29 -0400

To: John Young <jya@pipeline.com>

From: Rich Meislin <meislin@nytimes.com>

Subject: Re: CIA Iran Report

Dear Mr. Young, Thank you for informing us about the problem with this document. We are removing it from our site until we can delete the names in a more secure fashion. The names were obscured because of our concern for possible retribution against the families of the people named in this report, and we would strongly urge you to respect that judgment.

Sincerely, Rich Meislin

CRIME    APR 2016

# In online blunder, Dallas police revealed names of people reporting sexual assaults

**Sarah Mervosh, Investigative Reporter**

# Data can be revealing, even without names.

In March 2014, the New York City Taxi & License Commission tweeted a "TAXI FACTS" infographic:



Chris Whong files a "Freedom of Information Law" request for all the data used to create the graphic.

# NYC TLC provided Chris Whong with all of the data

175 million trips:

| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | medallion | hack_license | vendor_id | pickup_datetime | payment_type | fare_amoun | surcharge | mta_tax | tip_amount | tolls_amoun | total_amount |
| 2 | 89D227B655E5C82AECF13C3F | BA96DE419E711691B944 | CMT | 1/1/13 15:11 | CSH | 6.5 | 0 | 0.5 | 0 | 0 | 7 |
| 3 | 0BD7C8F5BA12B88E0B67BED | 9FD8F69F0804BDB5549F | CMT | 1/6/13 0:18 | CSH | 6 | 0.5 | 0.5 | 0 | 0 | 7 |
| 4 | 0BD7C8F5BA12B88E0B67BED | 9FD8F69F0804BDB5549F | CMT | 1/5/13 18:49 | CSH | 5.5 | 1 | 0.5 | 0 | 0 | 7 |
| 5 | DFD2202EE08F7A8DC9A57B0 | 51EE87E3205C985EF843: | CMT | 1/7/13 23:54 | CSH | 5 | 0.5 | 0.5 | 0 | 0 | 6 |
| 6 | DFD2202EE08F7A8DC9A57B0 | 51EE87E3205C985EF843: | CMT | 1/7/13 23:25 | CSH | 9.5 | 0.5 | 0.5 | 0 | 0 | 10.5 |
| 7 | 20D9ECB2CA0767CF7A01564 | 598CCE5B9C1918568DEE | CMT | 1/7/13 15:27 | CSH | 9.5 | 0 | 0.5 | 0 | 0 | 10 |
| 8 | 496644932DF3932605C22C7S | 513189AD756FF14FE670 | CMT | 1/8/13 11:01 | CSH | 6 | 0 | 0.5 | 0 | 0 | 6.5 |
| 9 | 0B57B9633A2FECD3D3B1944 | CCD4367B417ED6634D9{ | CMT | 1/7/13 12:39 | CSH | 34 | 0 | 0.5 | 0 | 4.8 | 39.3 |
| 10 | 2C0E91FF20A856C891483ED6 | 1DA2F6543A62B8ED9347 | CMT | 1/7/13 18:15 | CSH | 5.5 | 1 | 0.5 | 0 | 0 | 7 |

Every trip:

- Pickup date, time & GPS
- Drop-off date, time & GPS
- Fare & tip
- Encoded medallion number

https://en.wikipedia.org/wiki/Taxicabs_of_New_York_City

# With this data, you can make a map of NYC Taxi Service



Map of NYC, Plotted Using Locations Of All Yellow Taxi Dropoff
From January 2015 to June 2015

By Max Woolf — minimaxir.com     Made using R and ggplot2     Data via NYC TLC Trip Record Data 2015

http://minimaxir.com/2015/08/nyc-map/

# Compare taxi prices and Uber prices:

# Each taxi has a pseudonym, which allows taxi rides to be linked.

# Oops. The taxi medallion numbers were not properly de-identified.

| Pseudonym | Taxi Medallion |
|---|---|
| 0f76c35d4a069e0fe76b21d28f009639 | 5C27 |
| be9f314926dd314b36496d926e42f4db | 5C28 |
| 9ee993809f648d39d24f5ba8f862d7f1 | 5C29 |
| 23f7e8636fb9099822aa381054d215d4 | 5C30 |

The pseudonyms looked suspicious to Anthony Tockar, an intern at Neustar Research.

Tockar realized that the pseudonyms were MD5 hashes

MD5("5C28") = be9f314926dd314b36496d926e42f4db

# MD5 can't be reversed, but it's possible to do a "brute force search" on all possible values.

Anthony Tockar identified the medallion number the records.

He searched for photos in flickr that showed movie stars at taxis w



Riding with the Stars: Passenger Privacy in the NYC Taxicab Dataset

SEPTEMBER 15, 2014 BY ATOCKAR    56 COMMENTS

"5C27"

A journalist at Gawker identified 9 other cab rides.

# May 18, 1996: Massachusetts Governor William Weld Collapses at Bentley College Commencement



All·Politics
CNN TIME
News

## Massachusetts Governor Doing Well After Collapse

WALTHAM, Massachusetts (CNN, May 18) -- Gov. William Weld collapsed during a graduation ceremony at Bentley College, but doctors said he was doing well.

The governor was taken to Deaconess-Waltham Hospital, where he was undergoing a battery of tests, according to Bentley College spokeswoman Katherine Blake. Weld will remain in the hospital overnight for observation, she said.

Doctors said they performed an electrocardiogram, chest X-ray and blood tests, but found no immediate cause for concern.

"With all this testing we have done, nothing acute is showing," said Dr. Rifat Dweik.

"Right now, it looks like maybe the flu," said Pam Jonah, one of Weld's press aides.

Weld was receiving an honorary doctorate of law at 11 a.m. EDT when he was stricken, according to Blake.

# In 1997, MIT Graduate Student Latanya Sweeney decided to search for William Weld's medical records in the GIC data.

Sweeney obtains GIC dataset and looks for Weld's data.

- She knew that Weld lived in Cambridge, MA.
- Sweeney purchased Cambridge voter rolls for $20.
- Six people had the same birthday (July 31, 1945)
- Three were men
- One person had the same ZIP code.



02138

# "Linkage Attack"
# Matching records using quasi-identifiers

- Weld's records were uniquely identified.

- Sweeney estimated 87% of US population
  were uniquely identified by birthday, sex & ZIP



"Quasi-Identifiers"
or
"Indirect-Identifiers"

"Sensitive Data"

Hospital admission info

Birthday
Sex
ZIP Code

Name
Address
Phone
SSN

"Direct"
or
"Explicit"
identifiers

# Sweeney invented K-Anonymity
# A model for de-identifying structured data.

A dataset that you would like to release:

| Name | Race | Birthdate | Sex | Zip | Medication | Diagnosis |
|------|------|-----------|-----|-----|------------|-----------|
| Alice | Black | 9/20/65 | M | 37203 | M1 | Gastric Ulcer |
| Bob | Black | 2/14/65 | M | 37203 | M1 | Gastric Ulcer |
| Candice | Black | 10/23/65 | F | 37215 | M1 | Gastritis |
| Dan | Black | 8/24/65 | F | 37215 | M2 | Gastritis |
| Eliza | Black | 11/7/64 | F | 37215 | M2 | Gastritis |
| Felix | Black | 12/1/64 | F | 37215 | M2 | Stomach Cancer |
| Gazelle | White | 10/23/64 | M | 37215 | M3 | Flu |
| Harry | White | 3/15/64 | F | 37217 | M3 | Flu |
| Irene | White | 8/13/64 | M | 37217 | M3 | Flu |
| Jack | White | 5/5/64 | M | 37217 | M4 | Pneumonia |
| Kelly | White | 2/13/67 | M | 37215 | M4 | Pneumonia |
| Lenny | White | 3/21/67 | M | 37215 | M4 | Flu |

# First the identifiers are removed

| Name | Quasi Identifiers | | | | Medication | Diagnosis |
|------|------|------|------|------|------------|-----------|
| | Black | 9/20/65 | M | 37203 | M1 | Gastric Ulcer |
| | Black | 2/14/65 | M | 37203 | M1 | Gastric Ulcer |
| | Black | 10/23/65 | F | 37215 | M1 | Gastritis |
| | Black | 8/24/65 | F | 37215 | M2 | Gastritis |
| | Black | 11/7/64 | F | 37215 | M2 | Gastritis |
| | Black | 12/1/64 | F | 37215 | M2 | Stomach Cancer |
| | White | 10/23/64 | M | 37215 | M3 | Flu |
| | White | 3/15/64 | F | 37217 | M3 | Flu |
| | White | 8/13/64 | M | 37217 | M3 | Flu |
| | White | 5/5/64 | M | 37217 | M4 | Pneumonia |
| | White | 2/13/67 | M | 37215 | M4 | Pneumonia |
| | White | 3/21/67 | M | 37215 | M4 | Flu |

# A dataset is "k-anonymous" if every record is in a set of at least k indistinguishable individuals

| Race | Birthdate | Sex | Zip | Medication | Diagnosis |
|------|-----------|-----|-----|------------|-----------|
| Black | 65 | M | 37203 | M1 | Gastric Ulcer |
| Black | 65 | M | 37203 | M1 | Gastric Ulcer |
| Black | 65 | F | 37215 | M1 | Gastritis |
| Black | 65 | F | 37215 | M2 | Gastritis |
| Black | 64 | F | 37215 | M2 | Gastritis |
| Black | 64 | F | 37215 | M2 | Stomach Cancer |
| White | 64 | M | 3721- | M3 | Flu |
| White | 64 | - | 37217 | M3 | Flu |
| White | 64 | M | 3721- | M3 | Flu |
| White | 64 | - | 37217 | M4 | Pneumonia |
| White | 67 | M | 37215 | M4 | Pneumonia |
| White | 67 | M | 37215 | M4 | Flu |

# Attribute disclosure:
## We know the Black / 65 / M had a Gastric Ulcer.

| Race | Birthdate | Sex | Zip | Medication | Diagnosis |
|------|-----------|-----|-----|------------|-----------|
| Black | 65 | M | 37203 | M1 | Gastric Ulcer |
| Black | 65 | M | 37203 | M1 | Gastric Ulcer |
| Black | 65 | F | 37215 | M1 | Gastritis |
| Black | 65 | F | 37215 | M2 | Gastritis |
| Black | 64 | F | 37215 | M2 | Gastritis |
| Black | 64 | F | 37215 | M2 | Stomach Cancer |
| White | 64 | M | 3721- | M3 | Flu |
| White | 64 | - | 37217 | M3 | Flu |
| White | 64 | M | 3721- | M3 | Flu |
| White | 64 | - | 37217 | M4 | Pneumonia |
| White | 67 | M | 37215 | M4 | Pneumonia |
| White | 67 | M | 37215 | M4 | Flu |

# De-identification caveats — what can go wrong

Mistakes happen:

- Metadata may contain identifiers.
- Direct identifiers can be missed.
- Hard to determine what's a quasi-identifier.

Even worse:

- k-anonymity and I-diversity can significantly damage data quality.
- There is no mathematical proof that k-anonymity actually protects privacy.

Netflix Awards $1 Million Prize and Starts a New Contest

BY STEVE LOHR    SEPTEMBER 21, 2009 10:15 AM

Jason Kempin/Getty Images Netflix prize winners, from left: Yehuda Koren, Martin Chabbert, Martin Piotte, Michael Jahrer, Andreas Toscher, Chris Volinsky and Robert Bell.

**All data are potentially identifying.**

# The Netflix Challenge (2008-2009)

Netflix published movie data for ~450,000 subscribers:

- Pseudonymized username
- Information on movies watched:

*Movie Title*

*Date watched*

*Rating*

Challenge: Improve Netflix recommendation algorithm

Unintentional Challenge: Identify Netflix subscribers!

# Re-identifying the Netflix Challenge Victims



"Sensitive Data"

Other Movies Watched & Movie Rankings

Movies Watched & Movie Rankings

IMDB username

"Direct" or "Explicit" identifiers

Netflix Provided Data

United States Census Bureau

U.S. Department of Commerce
Economics and Statistics Administration
U.S. CENSUS BUREAU
census.gov

**Figure 4.** Adversary knows exact ratings and approximate dates.

**Figure 8.** Adversary knows exact ratings but does not know dates at all.

**Figure 9.** Effect of knowing less popular movies rated by victim. Adversary knows approximate ratings ($\pm 1$) and dates (14-day error).

# Netflix Settles Privacy Lawsuit, Cancels Prize Sequel



**The Firewall**
*the world of security* **FULL BIO** ∨

Opinions expressed by Forbes Contributors are their own.

**Taylor Buley**, Contributor

On Friday, Netflix announced on its corporate blog that it has settled a lawsuit related to its Netflix Prize, a $1 million contest that challenged machine learning experts to use Netflix's data to produce better recommendations than the movie giant could serve up themselves.

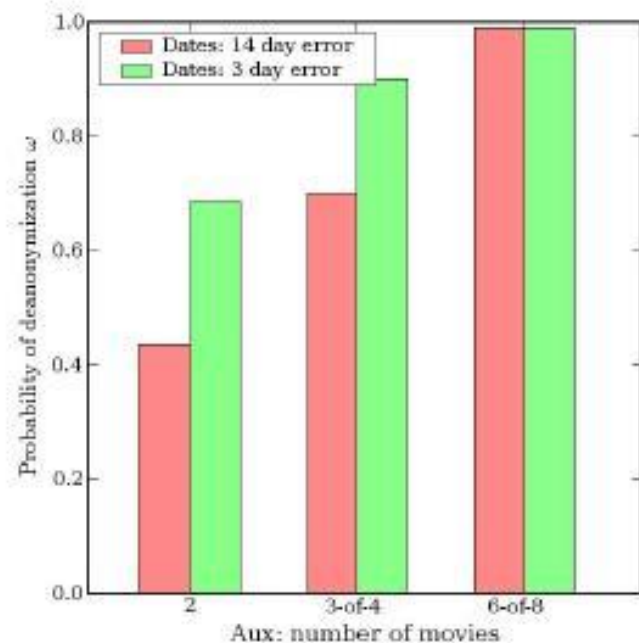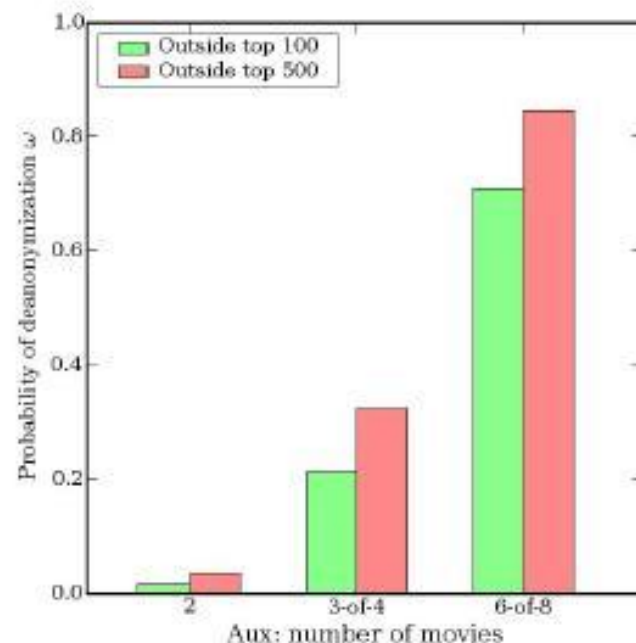The lawsuit called attention to academic research that suggests that Netflix indirectly exposed the movie preferences of its users by publishing anonymized customer data. In the suit, plaintiff Paul Navarro and others sought an injunction preventing Netflix from going through the so-called "Netflix Prize II," a follow-up challenge that Netflix promised would offer up even more personal data such as genders and zipcodes.

"Netflix is not going to pursue a sequel to the Netflix Prize," says spokesman Steve Swasey. "We looked at this, we heard some dissension and so we've settled it, resolved the issues and are moving on."

**United States Census** Bureau | **U.S. Department of Commerce**
Economics and Statistics Administration
U.S. CENSUS BUREAU
census.gov

**Differential Privacy: The Big Idea**

# Differential privacy is a new approach for assuring privacy in the release of statistical data.

Sensitive dataset

Ad hoc Rules

Privacy-Preserving Data Release

Privacy Parameters

Methods that implement the privacy definition

Formal Privacy Definition

Sensitive dataset

**Based on hope and assumptions.**

1. Data are identify, quasi-identifying, or not-identifying
2. Future data sets will not be released that can be linked with previously released data
3. Adversaries have limited resources to pursue re-identification attacks

**Based on math.**

# In traditional data publications, there are many ways that the contributions of an individual can leak out

**January**

| Name | Affect | Grade |
|---|---|---|
| Alex | Sad | 30 |
| Bobbie | Sad | 50 |
| Casey | Happy | 80 |
| Harper | Happy | 100 |

Statistical Tabulation →

Students: 4
Percent Happy: 50%
Average Grade: 65

**February**

| Name | Affect | Grade |
|---|---|---|
| Alex | Sad | 30 |
| Bobbie | Sad | 50 |
| Casey | Happy | 80 |
| Emerson | Sad | 90 |
| Harper | Happy | 100 |

Statistical Tabulation →

Students: 5
Percent Happy: 40%
Average Grade: 70

It's pretty easy to determine that the new kid is sad and has a 90.

**Differential privacy's core idea:**
**Create uncertainty regarding the presence any person in the dataset.**

# Noise is added to mask an individual's contribution

### January

| Name | Affect | Grade |
|---|---|---|
| Alex | Sad | 30 |
| Bobbie | Sad | 50 |
| Casey | Happy | 80 |
| Harper | Happy | 100 |

Statistical Tabulation + noise →

Students: 4
Percent Happy: 45%
Average Grade: 50

### February

| Name | Affect | Grade |
|---|---|---|
| Alex | Sad | 30 |
| Bobbie | Sad | 50 |
| Casey | Happy | 80 |
| Emerson | Sad | 90 |
| Harper | Happy | 100 |

Statistical Tabulation + noise →

Students: 5
Percent Happy: 60%
Average Grade: 75

# If we ran the statistics different times, we would get different results

January

| Name | Affect | Grade |
|------|--------|-------|
| Alex | Sad | 30 |
| Bobbie | Sad | 50 |
| Casey | Happy | 80 |
| Harper | Happy | 100 |

Statistical Tabulation + noise →

Students: 4
Percent Happy: 45%
Average Grade: 50

January

| Name | Affect | Grade |
|------|--------|-------|
| Alex | Sad | 30 |
| Bobbie | Sad | 50 |
| Casey | Happy | 80 |
| Harper | Happy | 100 |

Statistical Tabulation + noise →

Students: 4
Percent Happy: 55%
Average Grade: 75

January

| Name | Affect | Grade |
|------|--------|-------|
| Alex | Sad | 30 |
| Bobbie | Sad | 50 |
| Casey | Happy | 80 |
| Harper | Happy | 100 |

Statistical Tabulation + noise →

Students: 4
Percent Happy: 51%
Average Grade: 60

In this example, a *policy decision* requires that the number of students be accurately reported.

# Data users understand that noise has been added.

| Name | Affect | Grade |
|------|--------|-------|
| Alex | Sad | 30 |
| Bobbie | Sad | 50 |
| Casey | Happy | 80 |
| Harper | Happy | 100 |

January

Statistical Tabulation + noise →

Students: 3
Percent Happy: 40%
Average Grade: 50

| Name | Affect | Grade |
|------|--------|-------|
| Alex | Sad | 30 |
| Bobbie | Sad | 50 |
| Casey | Happy | 80 |
| Harper | Happy | 100 |

January

Statistical Tabulation + noise →

Students: 6
Percent Happy: 45%
Average Grade: 45

| Name | Affect | Grade |
|------|--------|-------|
| Alex | Sad | 30 |
| Bobbie | Sad | 50 |
| Casey | Happy | 80 |
| Harper | Happy | 100 |

January

Statistical Tabulation + noise →

Students: 5
Percent Happy: 51%
Average Grade: 60

In this example, a *policy decision* requires that the exact number of students in the class be confidential.

# How much noise do we add?
## That is a policy decision



epsilon vs. accuracy

Differential privacy uses the parameter ε (epsilon) to describe the privacy/accuracy tradeoff.

ε = 0  — No accuracy, full privacy
ε = ∞ — No privacy, full accuracy

# Noise can be added in two places:
# 1) When data are collected.   2) When statistics are produced.

Input noise infusion:

| Name | Affect | Grade |
|------|--------|-------|
| Alex | Sad + NOISE | 30 + NOISE |
| Bobbie | Sad + NOISE | 50 + NOISE |
| Casey | Happy + NOISE | 80 + NOISE |
| Harper | Happy + NOISE | 100 + NOISE |

Statistical Tabulation →

Students: 4
Percent Happy: 30..70
Average Grade: 50..80

Advantages:
- Tabulator need not be trusted.
- More statistics do not pose additional privacy threats.

Output noise infusion:

| Name | Affect | Grade |
|------|--------|-------|
| Alex | Sad | 30 |
| Bobbie | Sad | 50 |
| Casey | Happy | 80 |
| Harper | Happy | 100 |

Statistical Tabulation →

Students: 4
Percent Happy: 40..60
Average Grade: 60..70

Advantages:
- More accurate for the same level of privacy
- Allows uses of confidential data that do not involve publication.

# Other choices for policy makers
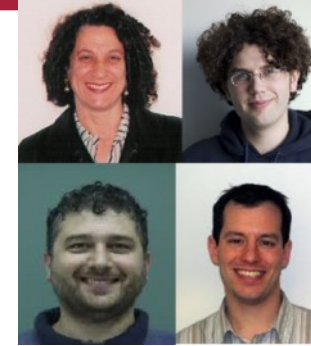
Where should the accuracy be spent?

What values should be reported exactly (with no privacy)

What are the possible bounds (sensitivity) of a person's data?
 e.g. If reporting average student age, can students be 5..18 or 5..115?

How do we convey privacy guarantees to public?

# Differential privacy was invented in 2006 by Dwork, McSherry, Nissim and Smith

Differential privacy is just 12 years old.



Today's public key cryptography was invented in 1976-1978

Remember public key cryptography in 1990?

- No standardized implementations. No SSL/TLS. No S/MIME or PGP.
- Very few people knew how to build systems that used crypto.

# In Summary

**Communications security**: Be careful when you get data form your sources.

**Storage security**: Be careful where you store data; use two-factor security.

**Publication security**: Be careful when you publish. Remember that data can be reverse-engineered if you do not take appropriate measures.

Questions?

Email: simson.l.garfinkel@census.gov