# Homomorphic Encryption in R

Alexander C. Mueller PhD
CEO and Founder of Capnion

May 9, 2019

# Who is the speaker?

30 second resume:

- an ancient metro Saint Louis townie
- grew up in University City
- University City High School
- B.A. Washington University (econ and math)
- Ph.D. University of Michigan (math)
- private sector data science
- founded data privacy company Capnion

# The Counter-Earth



Orbiting exactly opposite the Earth you are familiar with is
another Earth, alike in many ways but also different...

# Agenda

What and Why?

- What is homomorphic encryption?
- What's out there?
- Who cares? Use cases.

Examples and Application in R

- Basic keygen, encryption and decryption
- Easy in R, harder when R isn't R
- Wrinkles and how to flatten them
- An example classification problem
- Private models, private data

# Back to Basics



If only I could find someone with some **basic decency** and respect for other **people's privacy** to help me out.

I am not asking for anything too crazy, just someone competent in **basic arithmetic** to help with some **bookkeeping**.

Is that person you? Would you **care to volunteer**?
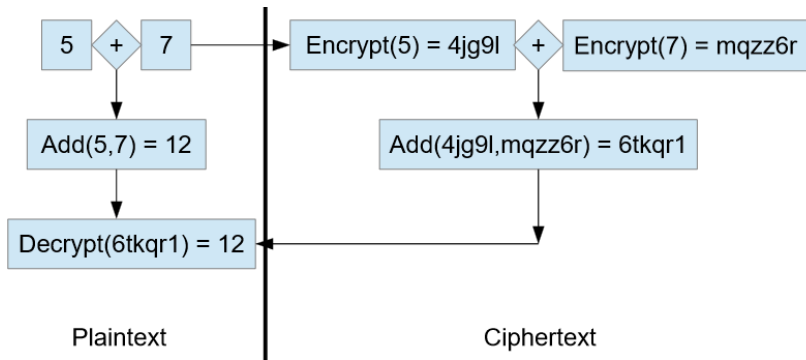
# Assignment: Basic Bookkeeping

Please go ahead and total the expenses by category.

| Expense | Amount | Starting Funds |
|---|---|---|
| fancy pens | XXXXXXXXX | XXXXXXXXX |
| private jet | XXXXXXXXX | **Ending Funds** |
| "business entertainment" | XXXXXXXXX | $30173 |
| attorneys | XXXXXXXXX | |

Easy, right? Thanks for handling this. It's so hard to find decent help these days. I need that jet deduction!

What's the problem here? Why is it hard to keep your private jet and business entertainment budget line items a secret?

# Homomorphic Encryption



Homomorphic encryption allows one to do computations on encrypted data and get the "correct" answer after decryption.

# History and Prophecy

Some landmark papers and people...

- (Gentry 2009) "A Fully Homomorphic Encryption Scheme"
- (Gennaro, Gentry, Parno, Raykova 2013) "Quadratic Span Programs and Succinct NIZKs without PCPs"

Craig Gentry won a MacArthur Fellowship (the "genius grant") for his 2009 PhD thesis work on fully homomorphic encryption (able to compute an arbitrary algorithm on ciphertext) and was widely interviewed afterwards.

He said at the time, almost 10 years ago, that it might be 10 years before the technology really started to catch on...

# What's Out There?

- HomomorphicEncryption R Package
  - intuitive interface in R
  - http://www.louisaslett.com/HomomorphicEncryption/
- HElib - fully homomorphic encryption software library
  - the algorithms, close to the metal in C++
  - https://github.com/shaih/HElib
- PySEAL Python Package
  - intuitive interface in Python
  - https://github.com/Lab41/PySEAL
- OpenMined - crowdsourced, private machine learning
  - https://github.com/OpenMined

# The HomomorphicEncryption Library

"For research purposes only"

- Created by academic statistician Louis Aslett [1]
- A wrapper for a C++ library, via Fan and Vercauteren [2]

Some notable limitations

- Exists for Mac and Linux but not Windows
- Pain around where R is really C under the hood
- (True of all HE:) More algebra, More problems

Homomorphic encryption is a new technology, and there is work still to be done bridging the space between basic artithmetic and scientific computing.

# Business Layer: Analytics Partners

Businesses (i.e. insurance companies) often hand data off to **analytics partners** (i.e. risk scorers) to obtain some insight on that data. For a variety of reasons, the primary data holder might want to keep their data secret and / or the analytics partner might want to keep their model secret.

Depending on the industry, there are many businesses that fit into this rough category...

- **Risk scoring in insurance**
- Raters of consumer credit
- Clearing houses for financial transaction
- Many forms of consulting

# Typical Data Flow

# Data Flow with Homomorphic Encryption

# R Code 1: Getting Started



```
bird@python-entity: ~                                              —  □  ×

> #load the library
> library("HomomorphicEncryption")
Loading required package: gmp

Attaching package: âgmpâ

The following objects are masked from âpackage:baseâ:

    %*%, apply, crossprod, matrix, tcrossprod

NOTE: this is an academic research implementation of homomorphic encryption sche
mes and no warranty of correctness or security is provided.

For citation information, type citation("HomomorphicEncryption").


Attaching package: âHomomorphicEncryptionâ

The following objects are masked from âpackage:baseâ:

    save, save.image, saveRDS

> #built on a C library
> #uses a particular cryptosystem
> #this is Fan and Vercauteren
> p <- pars("FandV")
> print(p)
Fan and Vercauteren parameters
Ï = xâ´â°â³â¶+1
q = 3402823669209384634633746074317682111456 (128-bit integer)
t = 32768
Î = 10384593717069655257060992658440192
Î = 16
Security level â 128-bits
Supports multiplicative depth of 3 with overwhelming probability (i.e. lower bou
nd, likely more possible)
>
```

# R Code 2: Key Generation



We end up with an asymmetric key pair with special powers.

# R Code 3: New Data Types



```
bird@python-entity: ~                        —    □    ×
> #encrypt a couple vectors
> c1 <- enc(k$pk, c(42,34))
> c2 <- enc(k$pk, c(7,5))
>
> #the data is encrypted
> print(c2)
Vector of 2 Fan and Vercauteren cipher texts
> #working with a special data type
> print(typeof(c2))
[1] "S4"
>
```

We have new data types for individual pieces of ciphertext and
also vectors thereof. They will do their best to pretend to be the
analogous base R data types but they are not.

# R Code 4: Operator Overloading



```
bird@python-entity: ~                    —    □    ×

> #can add, multiply, inner product
> #R syntax is the same
> cres1 <- c1 + c2
> cres2 <- c1 * c2
> cres3 <- c1 %*% c2
>
> #decrypt and get correct answers
> dec(k$sk, cres1)
[1] 49 39
> dec(k$sk, cres2)
[1] 294 170
> dec(k$sk, cres3)
[1] 464
>
```

R will add our encrypted numbers, and vectors thereof, using
the usual syntax. Any code written purely in R using this syntax
will run and get the "right" answer on encrypted data.

# R Code 4.5: What really just happened?



These polynomials are what we just "added" together. When we decrypted, whatever we did to those polynomials turned out to be the right thing in that we managed to compute addition of the plaintext correctly.

# R Code 5: Some Wrinkles
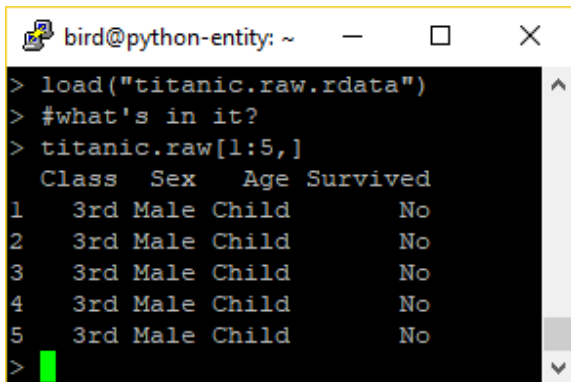


```
bird@python-entity: ~                          —    □    ×
> #vector like but not an atomic vector
> vectorLike <- enc(k$pk, c(42,34))
> #can get length for example
> length(vectorLike)
[1] 2
> #and slice as we would a vector
> dec(k$sk,vectorLike[1]+vectorLike[2])
[1] 76
> #this doesn't work
> enc(k$pk, 5.5)
Error in enc.Rcpp_FandV_pk(k$pk, 5.5) : Only integers can be encrypted.
>
```

The cipher **only works for integers**, which is a huge pain in the ass in practice but not so bad in principle.

You can **approximate by fractions**, and computations on fractions are just computations on integers as $\frac{a}{b} + \frac{c}{d} = \frac{ad+bc}{db}$.

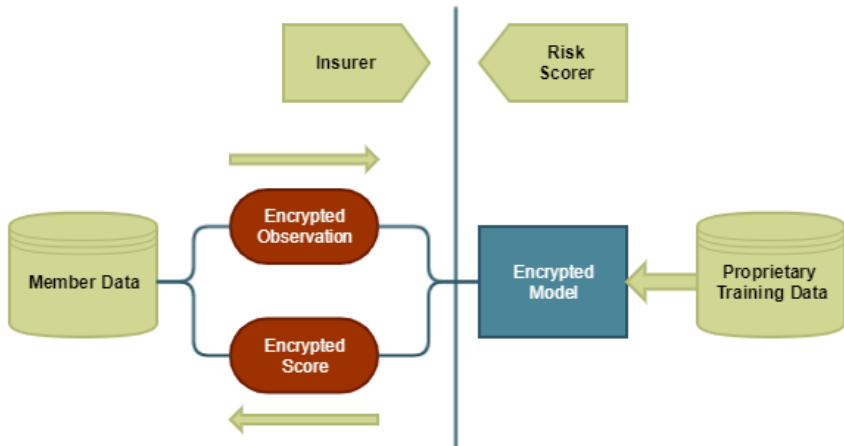This is not the level of abstraction we are used to, though.
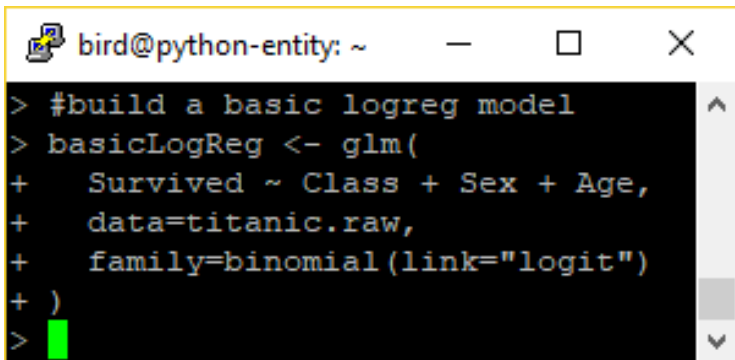
# R Code 6: A Familiar Classification Problem



```
bird@python-entity: ~              —    □    ×
> load("titanic.raw.rdata")
> #what's in it?
> titanic.raw[1:5,]
  Class  Sex   Age Survived
1  3rd Male Child       No
2  3rd Male Child       No
3  3rd Male Child       No
4  3rd Male Child       No
5  3rd Male Child       No
>
```

We'll train a logistic regression model on the Titanic dataset,
encrypt the data of that model, and compute survival odds
**blind to both the observations used and the results**.
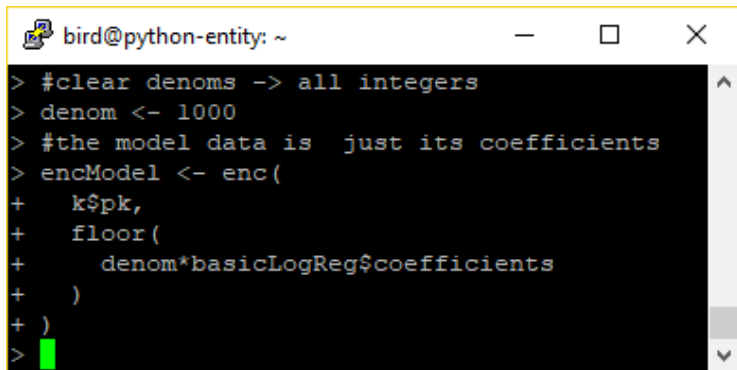
# Encrypted Model Data Flow

# R Code 7: Train on Plaintext

```
bird@python-entity: ~                    —    □    ×
> #build a basic logreg model
> basicLogReg <- glm(
+     Survived ~ Class + Sex + Age,
+     data=titanic.raw,
+     family=binomial(link="logit")
+ )
>
```

Here we imagine we had some in-house data we could use to train our model as usual, but maybe we have clients that don't want to share information about themselves.
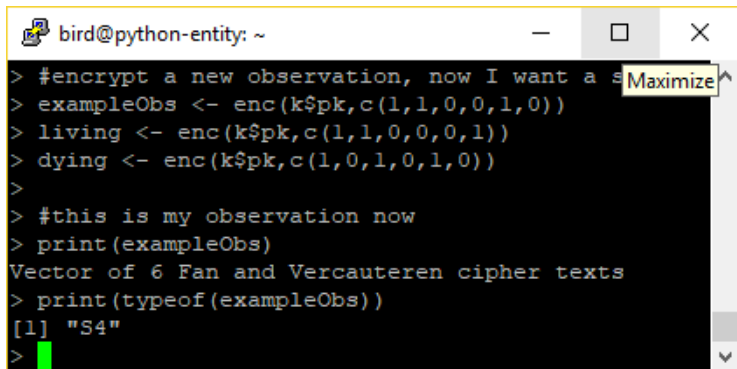
# R Code 8: Encrypt the Model



```
> #clear denoms -> all integers
> denom <- 1000
> #the model data is  just its coefficients
> encModel <- enc(
+    k$pk,
+    floor(
+      denom*basicLogReg$coefficients
+    )
+ )
>
```

A logistic regression model is determined entirely by its coefficients, which we can encrypt like anything else after clearing denominators (-ish) so we can work with integers.

# R Code 9: Encrypted Observations



```
bird@python-entity: ~                           —    □    ×
> #encrypt a new observation, now I want a s  Maximize ^
> exampleObs <- enc(k$pk,c(1,1,0,0,1,0))
> living <- enc(k$pk,c(1,1,0,0,0,1))
> dying <- enc(k$pk,c(1,0,1,0,1,0))
>
> #this is my observation now
> print(exampleObs)
Vector of 6 Fan and Vercauteren cipher texts
> print(typeof(exampleObs))
[1] "S4"
>
```

We'll reformat a bit as we are going to write our own prediction
function, and then we can take example observations and
encrypt them so the model owner need not see our business.

# R Code 10: Insights



```
bird@python-entity: ~                              —    □    ×
> plain_odds <- function(obs,coeffs,sk,denom){
+    innerProd <- dec(sk,obs %*% coeffs)
+    decLogOdds <- (1/denom)*innerProd
+    return(exp(decLogOdds)/(1+exp(decLogOdds)))
+ }
> print(plain_odds(exampleObs,encModel,k$sk,denom))
[1] 0.1982926
> print(plain_odds(living,encModel,k$sk,denom))
[1] 0.8894367
> print(plain_odds(dying,encModel,k$sk,denom))
[1] 0.1037719
>
```

We can now compute our model, defined by encrypted data, on encrypted observations to get an encrypted answer. When we decrypt, we see that we are getting good results.
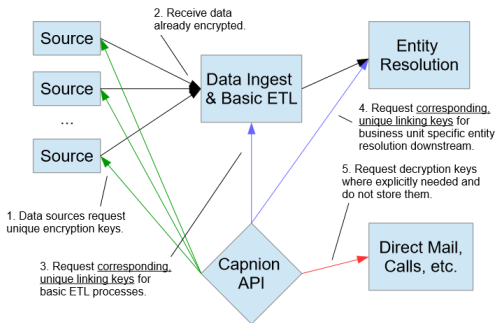
Let's say I have an encrypted string and an algorithm that takes a plaintext search string and returns the indices where that string appears.

$$\text{search("a")} = \{2, 8, 13\}$$

$$\text{search("b")} = \{1\}$$

$$\text{search("c")} = \{\}$$

and so on ...

It will never take me more than 26 tries before I can reconstruct **"bart skateboards"** as the plaintext.

# Security Subtleties 2



Generate keys that enable **particular operations on particular sets of data**, and specifically keys that give information on **how to link records but not on how to decrypt** or do other computations

# Questions

Any questions?

Feel free to contact me at acmueller@capnion.com

Slides are available at
https://github.com/capnion/ ...
random/blob/master/acm_rug_he_may19.pdf

# References

📄 L. J. M. Aslett, P. M. Esperança, and C. C. Holmes, "A review of homomorphic encryption and software tools for encrypted statistical machine learning," tech. rep., University of Oxford, 2015.

📄 J. Fan and F. Vercauteren, "Somewhat practical fully homomorphic encryption," tech. rep., 2012.