

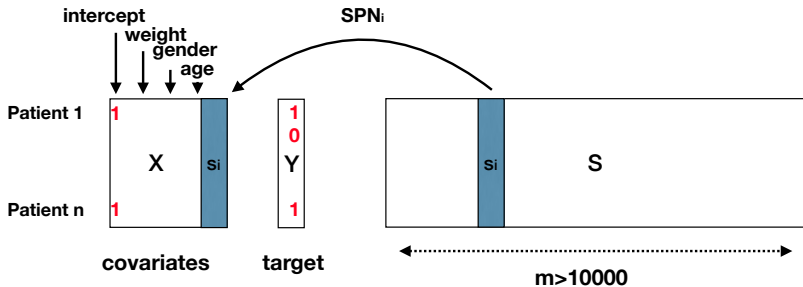
Privacy-preserving Efficient Subset of Features Selection for Regression Models

N. Gama, M. Georgieva



inpher.io

GWAS (find the best additional feature)



Question: Is the new feature important?

Naive method: compute $stat_i$ for each i ...

... that means compute more than 10000 logreg

Goal:

Develop a secure parallel outsourcing solution to compute Genome Wide Association Studies (GWAS) based on linear/logistic regression using **homomorphically encrypted** data.

Challenge (informally):

- Currently: a logreg model, 250 patients, with 3 or 4 physical features.
- Which new feature, among 10000 possible genes (SNP), would improve the model?

Semi-parallel approach:

- Don't do 10000 logregs...

Goal:

Develop a secure parallel outsourcing solution to compute Genome Wide Association Studies (GWAS) based on linear/logistic regression using **homomorphically encrypted** data.

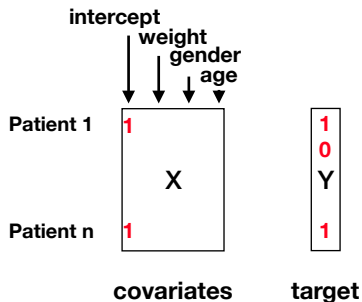
Challenge (informally):

- Currently: a logreg model, 250 patients, with 3 or 4 physical features.
- Which new feature, among 10000 possible genes (SNP), would improve the model?

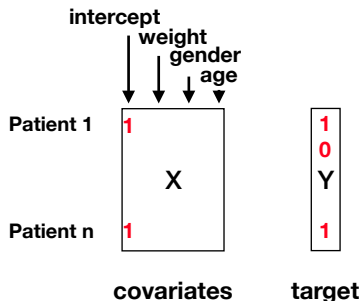
Semi-parallel approach

- Don't do 10000 logregs...

Logreg, IRLS, relevance of a feature



- Single Logistic regression:
 - Find θ s.t $Y = \text{sign}(X\theta)$
- IRLS:
 - Compute $\text{grad} = X^t(Y - p)$, with $p = \sigma(X\theta)$
 - Compute $\text{Hessian} = X^t \text{diag}(p(1 - p))X$



Importance of the i^{th} feature:

- the i^{th} coeff is big: θ_i (numerator)
- the i^{th} error term is small:
 $(Hess^{-1})_{i,i}$ (denominator)

stat= ratio

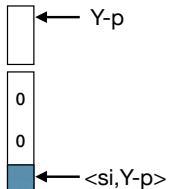
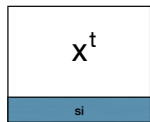
- Single Logistic regression:
 - Find θ s.t $Y = sign(X\theta)$
- IRLS:
 - Compute $grad = X^t(Y - p)$, with $p = \sigma(X\theta)$
 - Compute $Hessian = X^t diag(p(1 - p))X$

Semi-parallel GWAS (high level idea)

Semi-parallel GWAS (optimized)

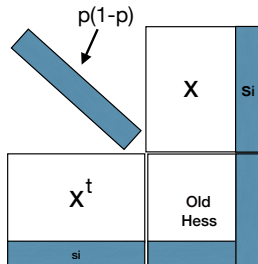
- 1 Do $\text{logreg}(X, y)$ without S
- 2 Once model is converged, add s_i

- Gradient:



They can be batch-computed: $(Y-p)^t S$

- Hessian:



FHE

- Long term storage
- Unique Cloud
- Slower and consumes more memory

MPC

- Faster than FHE
- More accuracy
- All data owner must participate

Fixed points versus Floating point

Floating point:

- $x = m \cdot 2^\tau$, with $m \in 2^{-\rho} \cdot \mathbb{Z}$ and $\frac{1}{2} \leq |m| < 1$
- $\tau = \lceil \log_2(x) \rceil$ data dependent and **not public (not FHE-friendly)**
- **The exponent is always in sync with the data**
ex: $(1.23 \cdot 10^{-4}) * (7.24 \cdot 10^{-4}) = (8.90 \cdot 10^{-8})$

Fixed point:

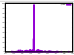

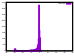
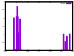
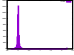

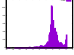
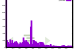
- $x = m \cdot 2^\tau$, with $m \in 2^{-\rho} \cdot \mathbb{Z}$ and $0 \leq |m| < 1$,
- τ **is public, thus FHE-friendly**
- **Risk of overflow** (τ too small)
- **Risk of underflow** (τ too large)
ex: $(0.000123 \cdot 10^0) * (0.000724 \cdot 10^0) = (0.000000 \cdot 10^0)$

Plaintext parameters:

- $\rho \in \mathbb{N}$: bits of precision of the plaintext (≈ 15 bits)
- $\tau \in \mathbb{Z}$: slot exponent (order of magnitude of the complex values in each slot)

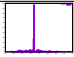

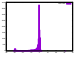
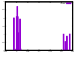
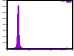

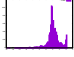

Choice of slot exponent

The slot exponent τ that defines the plaintext interval must be carefully estimated.

variable	avg	stdev	min	max	dist
p	0.440816	0.0975715	0.176397	0.853487	
w	0.236977	0.0201871	0.125047	0.25	
z_i^*	-3.33092	7.36068	-30.9426	31.2008	
G	0.0577846	0.0953495	-0.011997	0.236977	
A	0.0621965	0.301255	-0.317312	2.236	
$(s_i^*)^2$	2.44243	4.11085	0.111961	14.5044	
$\log(\text{stat}_i)$	0.200039	1.84459	-13.7207	4.36158	
p -value	0.310218	0.24083	0	0.999163	

Choice of slot exponent

The slot exponent τ that defines the plaintext interval must be carefully estimated.

variable	avg	stdev	min	max	dist
p	0.440816	0.0975715	0.176397	0.853487	
w	0.236977	0.0201871	0.125047	0.25	
z_i^*	-3.33092	7.36068	-30.9426	31.2008	
G	0.0577846	0.0953495	-0.011997	0.236977	
A	0.0621965	0.301255	-0.317312	2.236	
$(s_i^*)^2$	2.44243	4.11085	0.111961	14.5044	
$\log(\text{stat}_i)$	0.200039	1.84459	-13.7207	4.36158	
p -value	0.310218	0.24083	0	0.999163	



Not stable

Increase the precision of the algorithm, but that implies bigger parameters.



Stable

Use stable computation with negative feedback
(e.g. gradient descent)
Smaller parameters

FHE parameters:

- $L \in \mathbb{N}$: level exponent of the ciphertext ($\alpha = 2^{-(L+\rho)}$: noise rate)
- $N = f(\lambda, \alpha)$: key size, with λ the security parameter

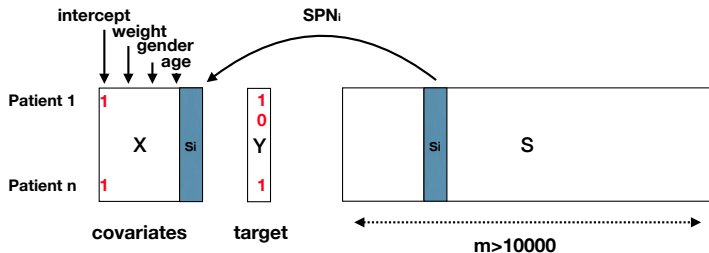
The `lwe-estimator` script was used to assert the security.
(conform to HE security standardization white paper)

FHE parameters:

- $L \in \mathbb{N}$: level exponent of the ciphertext ($\alpha = 2^{-(L+\rho)}$: noise rate)
- $N = f(\lambda, \alpha)$: key size, with λ the security parameter

The `lwe-estimator` script was used to assert the security.
(conform to HE security standardization white paper)

Plaintext algorithm in FHE solution



Input:

- $X \in \mathcal{M}_{n,k+1}(\mathbb{R})$ input matrix
- $y \in \mathbb{B}^n$ binary vector
- $S \in \mathcal{M}_{n,m}(\mathbb{R})$ assumed binary

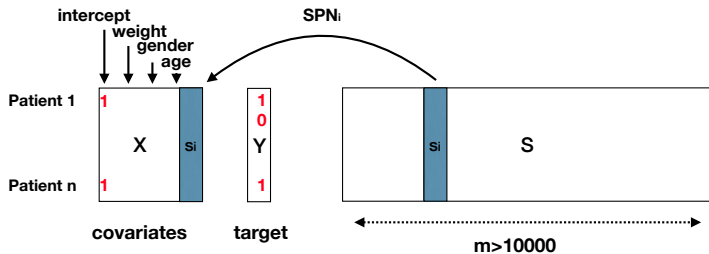
Output:

- $stat \in \mathbb{R}^m$
with $stat_i = \frac{z^*_i}{\sqrt{s^*_i}}$

Key points of our solution:

- Make plaintext algorithm FHE friendly
- Use hybrid homomorphic encryption

Plaintext algorithm in FHE solution



Input:

- $X \in \mathcal{M}_{n,k+1}(\mathbb{R})$ input matrix
- $y \in \mathbb{B}^n$ binary vector
- $S \in \mathcal{M}_{n,m}(\mathbb{R})$ assumed binary

Output:

- $stat \in \mathbb{R}^m$
with $stat_i = \frac{z^*_i}{\sqrt{s^*_i}}$

Key points of our solution:

- Make plaintext algorithm FHE friendly
- Use hybrid homomorphic encryption

Make the plaintext algorithm FHE friendly

- Find simple geometric equivalents of the formula
- Find approximation with lower multiplicative depth
- Replace feature scaling of X with orthogonalization

Algorithm 2 Plaintext algorithm

- 1: $\beta^{(0)} = (0, \dots, 0)$
 - 2: **for** $t = 1$ **to** iters **do**
 - 3: $\beta^{(t)} \leftarrow \beta^{(t-1)} + \text{step} \cdot X^T \cdot (\mathbf{y} - \sigma(X\beta^{(t-1)}))$
 - 4: **end for** ▷ logreg

 - 5: $\mathbf{p} \leftarrow \sigma(X\beta^{(\text{iters})})$
 - 6: $\mathbf{z}^* \leftarrow (\mathbf{y} - \mathbf{p})^T \cdot S$ ▷ numerator
 - 7: $W \leftarrow \text{diag}(p * (1 - p))$
 - 8: $G \leftarrow X^T \cdot W \cdot X \approx \frac{1}{4} * Id$ (assumed that X orthogonal)
 - 9: $A \leftarrow X^T \cdot W \cdot S$
 - 10: $\mathbf{s}^{*2} = \text{colsums}(W \cdot (S \odot S)) - \text{colsums}(A \odot G^{-1}A)$
 - 11: $(\approx A[0] * \sqrt{n} - 4 * \text{colsums}(A \odot A))$ ▷ denominator

 - 12: $\mathbf{r}_i = [2 \cdot \log(|\mathbf{z}^*_i|) - \log(|\mathbf{s}^{*2}_i|)]$ for each $i \in [1, m]$ ▷ log of stat
-

$$\text{stat}_i = \frac{\mathbf{z}^*_i}{\sqrt{\mathbf{s}^{*2}_i}} = \frac{1}{2} \exp(\mathbf{r}_i)$$

$$\text{p-value}_i = \text{p-Norm}(\text{stat}_i)$$

Algorithm

for loops

(better with fast bootstrapping)

- 1: $\beta^{(0)}$
- 2: **for** $t = 1$ **to** iters **do**
- 3: $\beta^{(t)} \leftarrow \beta^{(t-1)} + \text{step} \cdot X^T \cdot (\mathbf{y} \rightarrow \sigma(X\beta^{(t-1)}))$
- 4: **end** ▷ logreg
- 5: \mathbf{p} ▷ numerator
(Approx numbers, or Lookup tables)
- 6: \mathbf{z}^*
- 7: $W \leftarrow \text{diag}(p * (1 - p))$
- 8: $G \leftarrow X^T \cdot W \cdot X \approx \frac{1}{4} * Id$ (assumed that X orthogonal)
- 9: $A \leftarrow X^T \cdot W \cdot S$
- 10: $\mathbf{s}^{*2} = \text{colsums}(W \cdot (S \odot S)) - \text{colsums}(A \odot G^{-1}A)$
- 11: $(\approx A[0] * \sqrt{n} - 4 * \text{colsums}(A \odot A))$ ▷ denominator
- 12: $\mathbf{r}_i = [2 \cdot \log(|\mathbf{z}^*_i|) - \log(|\mathbf{s}^{*2}_i|)]$ for each $i \in [1, m]$ ▷ log of stat

$$\text{stat}_i = \frac{\mathbf{z}^*_i}{\sqrt{\mathbf{s}^{*2}_i}} = \frac{1}{2} \exp(\mathbf{r}_i)$$

$$\text{p-value}_i = \text{p-Norm}(\text{stat}_i)$$

Algorithm

for loops
(better with fast bootstrapping)

- 1: $\beta^{(0)}$
- 2: **for** $t = 1$ **to** iters **do**
- 3: $\beta^{(t)} \leftarrow \beta^{(t-1)} + \text{step} \cdot X^T \cdot (\mathbf{y} \rightarrow \sigma(X\beta^{(t-1)}))$
- 4: **end**
- 5: \mathbf{p}
- 6: \mathbf{z}^*
- 7: $W \leftarrow \text{diag}(\mathbf{p})$
- 8: $G \leftarrow X^T \cdot W$
- 9: $A \leftarrow X^T \cdot W \cdot S$
- 10: $\mathbf{s}^{*2} = \text{colsums}(W \cdot (S \odot S)) - \text{colsums}(A \odot G^{-1}A)$
- 11: $(\approx A[0] * \sqrt{n} - 4 * \text{colsums}(A \odot A))$
- 12: $\mathbf{r}_i = [2 \cdot \log(|\mathbf{z}^*_i|) - \log(|\mathbf{s}^{*2}_i|)]$ for each $i \in [1, m]$

continuous non-polynomial functions
(Approx numbers, or Lookup tables)

individual non-linear operations in small dimension
(lookup tables)

multiplication with fresh ciphertexts
(better with TFHE's external product)

▷ numerator

▷ denominator

▷ log of stat

$$\text{stat}_i = \frac{\mathbf{z}^*_i}{\sqrt{\mathbf{s}^{*2}_i}} = \frac{1}{2} \exp(\mathbf{r}_i)$$

$$\text{p-value}_i = \text{p-Norm}(\text{stat}_i)$$

Algorithm in plaintext



Algorithm

for loops
(better with fast bootstrapping)

- 1: $\beta^{(0)}$
- 2: **for** $t = 1$ **to** iters **do**
- 3: $\beta^{(t)} \leftarrow \beta^{(t-1)} + \text{step} \cdot X^T \cdot (\mathbf{y} \rightarrow \sigma(X\beta^{(t-1)}))$
- 4: **end**
- 5: \mathbf{p}
- 6: \mathbf{z}^*
- 7: $W \leftarrow \text{diag}(\mathbf{p})$
- 8: $G \leftarrow X^T \cdot W$
- 9: $A \leftarrow X^T \cdot W \cdot S$
- 10: $\mathbf{s}^{*2} = \text{colsums}(W \cdot (S \odot S)) - \text{colsums}(A \odot G^{-1}A)$
- 11: $(\approx A[0] * \sqrt{n} - 4 * \text{colsums}(A \odot A))$
- 12: $\mathbf{r}_i = [2 \cdot \log(|\mathbf{z}^*_i|) - \log(|\mathbf{s}^{*2}_i|)]$ for each $i \in [1, m]$

continuous non-polynomial functions
(Approx numbers, or Lookup tables)

individual non-linear operations in small dimension
(lookup tables)

multiplication with fresh ciphertexts
(better with TFHE's external product)

very large dimension
(fully packed SIMD)

▷ numerator

$$\text{stat}_i = \frac{\mathbf{z}^*_i}{\sqrt{\mathbf{s}^{*2}_i}} = \frac{1}{2} \exp(\mathbf{r}_i)$$

$$\text{p-value}_i = \text{p-Norm}(\text{stat}_i)$$

continuous function batched on a large vector

Algorithm in plaintext

Algorithm

- 1: $\beta \leftarrow 0$
- 2: **for** $t = 1$ **to** iters **do**
- 3: $\beta^{(t)} \leftarrow \beta^{(t-1)} + \text{step} \cdot X^T \cdot (\mathbf{y} \rightarrow \sigma(X\beta^{(t-1)}))$
- 4: **end for**
- 5: $\mathbf{p} \leftarrow \text{colsums}(X \cdot \beta^{(t)})$
- 6: $\mathbf{z}^* \leftarrow \text{colsums}(X \cdot \mathbf{p})$
- 7: $W \leftarrow \text{diag}(\mathbf{z}^*)$
- 8: $G \leftarrow X^T \cdot W$
- 9: $A \leftarrow X^T \cdot W \cdot X$
- 10: $\mathbf{s}^{*2} = \text{colsums}(W \cdot (S \odot S)) - \text{colsums}(A \odot G^{-1}A)$
- 11: $(\approx A[0] * \sqrt{n} - 4 * \text{colsums}(A \odot A))$
- 12: $\mathbf{r}_i = [2 \cdot \log(|\mathbf{z}^*_i|) - \log(|\mathbf{s}^{*2}_i|)]$ for each $i \in [1, m]$

for loops
(better with fast bootstrapping)

continuous non-polynomial functions
(Approx numbers, or Lookup tables)

individual non-linear operations in small dimension
(lookup tables)

multiplication with fresh ciphertexts
(better with TFHE's external product)

▷ numerator

very large dimension
(fully packed SIMD)

$$\text{stat}_i = \frac{\mathbf{z}^*_i}{\sqrt{\mathbf{s}^{*2}_i}} = \frac{1}{2} \exp(\mathbf{r}_i)$$

$$\text{p-value}_i = \text{p-Norm}(\text{stat}_i)$$

continuous function batched on a large vector

Which fully homomorphic scheme should we choose?

Each library has its own strengths



Strengths of HE libraries

- BGV/Helib: SIMD finite field arithmetic
- B/FV, Seal: SIMD vector $\text{mod } t$
- HEAAN: SIMD fixed point arithmetic
- TFHE: single evaluation, boolean logic, comparison, threshold, complex circuits
- etc...

How to get all the benefits without the limitations?

Idea:

- Unified plaintext space over the Torus
- Switch between ciphertext representations
- Implement bridges between TFHE, B/FV and HEAAN

For this use-case

We use the switch between TFHE and HEAAN!

Idea:

- Unified plaintext space over the Torus
- Switch between ciphertext representations
- Implement bridges between TFHE, B/FV and HEAAN

For this use-case

We use the switch between TFHE and HEAAN!

- 1 Initial Logreg on matrix X and vector y
 - adapt lib TFHE + logreg
- 2 Mass Linear algebra computations
 - implement Chimera (version 2 of TFHE)
- 3 Batch Logarithm computation
 - adapt lib HEAAN

Steps	Timing (4 cores)	Timing (96 cores)	RAM
KeyGen	5.5 mins	2.0 mins	4.4 GB
Encryption	7.2 mins	1.3 mins	8.6 GB
Cloud Computation	3h06	10.2 mins	7.8 GB

- Input ciphertext: 5GB (enc X, y, S)
- Final ciphertext: 640KB (enc numerator + denominator)

Benchmarks (with new optimizations)



$k = 3, n = 250, m = 10000$

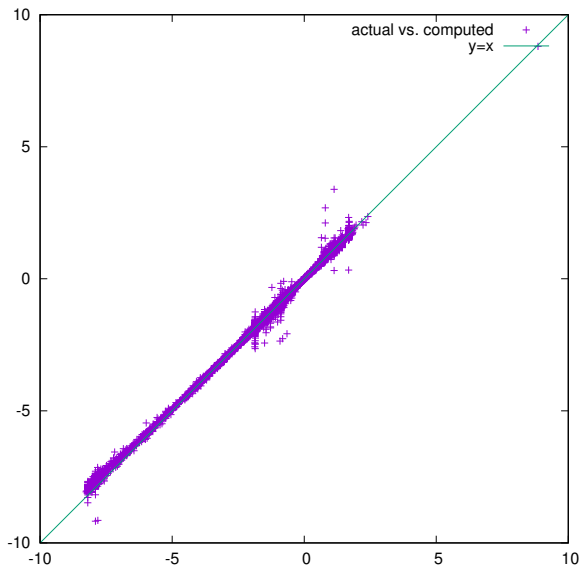
Steps	Timing (4 cores)	Timing (96 cores)	RAM
KeyGen	5.5 mins	2.0 mins	4.4 GB
Encryption	7.2 mins	1.3 mins	8.6 GB
Cloud Computation	35 mins	3 mins	7.8 GB

$k = 7, n = 250, m = 10000$

Steps	Timing (4 cores)	Timing (96 cores)	RAM
KeyGen	5.5 mins	2.0 mins	4.4 GB
Encryption	7.2 mins	1.3 mins	8.6 GB
Cloud Computation	41 mins	3.1 mins	7.8 GB

- initial ciphertext: 5GB (enc X, y, S)
- final ciphertext: 640KB (enc numerator + denominator)

Numerical Accuracy (FHE has noise)



Questions?

