# Future Perspectives of Privacy-Preserving Machine Learning

Rafael Dowsley

# Machine Learning: Powerful Tool

Enables to automatically learn from huge amounts of data and make sense of a complicated world.

Has the potential to vastly improve the quality of our daily lives:

- Diagnose patients

- Personalize search engines and recommendation platforms
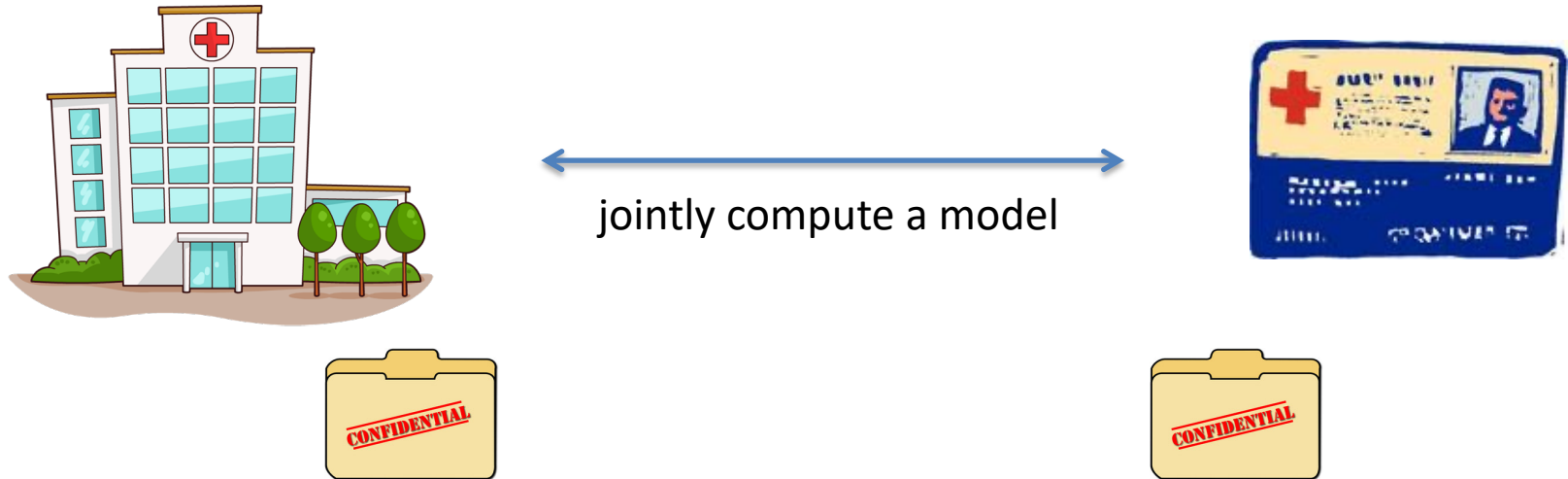
- Self-driving cars

- Optimize healthcare spending

- …

# Dilemma

Traditional machine learning methods assume that <span style="color:red">all existing data</span> is available.

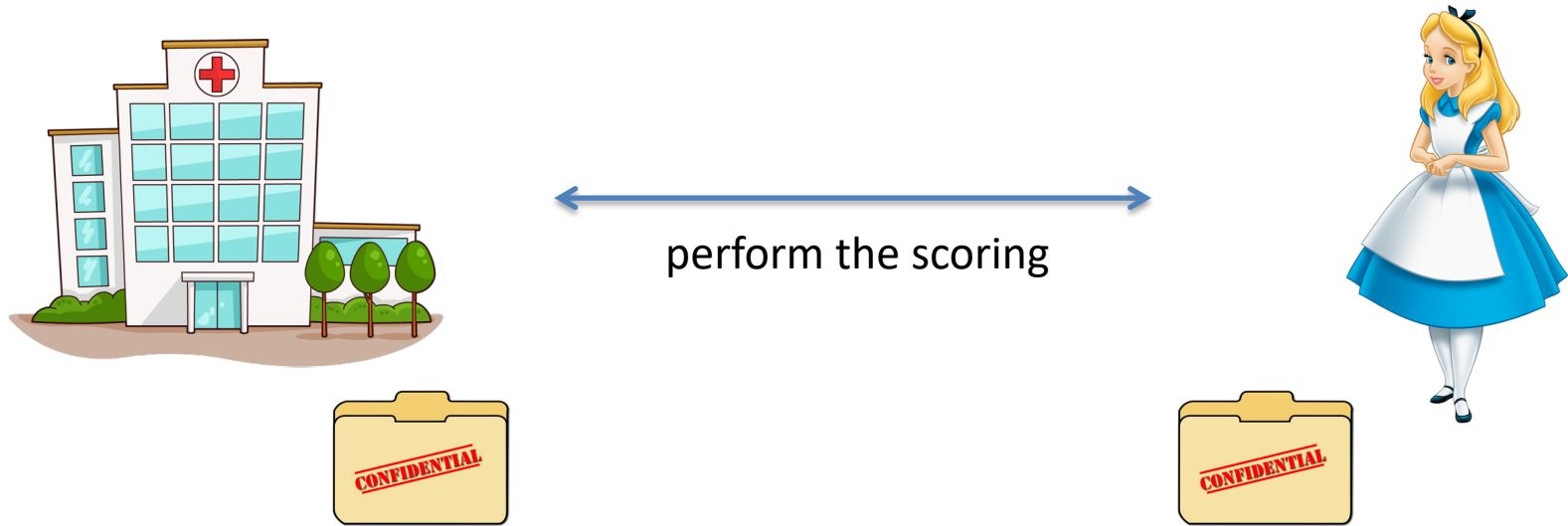The data needed to train the model can be owned by many different parties.

Also in the task of performing a scoring, there can also be privacy issues between the model owner and the user.

# Privacy-Preserving Learning



jointly compute a model

They cannot/do not want to share their data set.

# Privacy-Preserving Scoring

perform the scoring

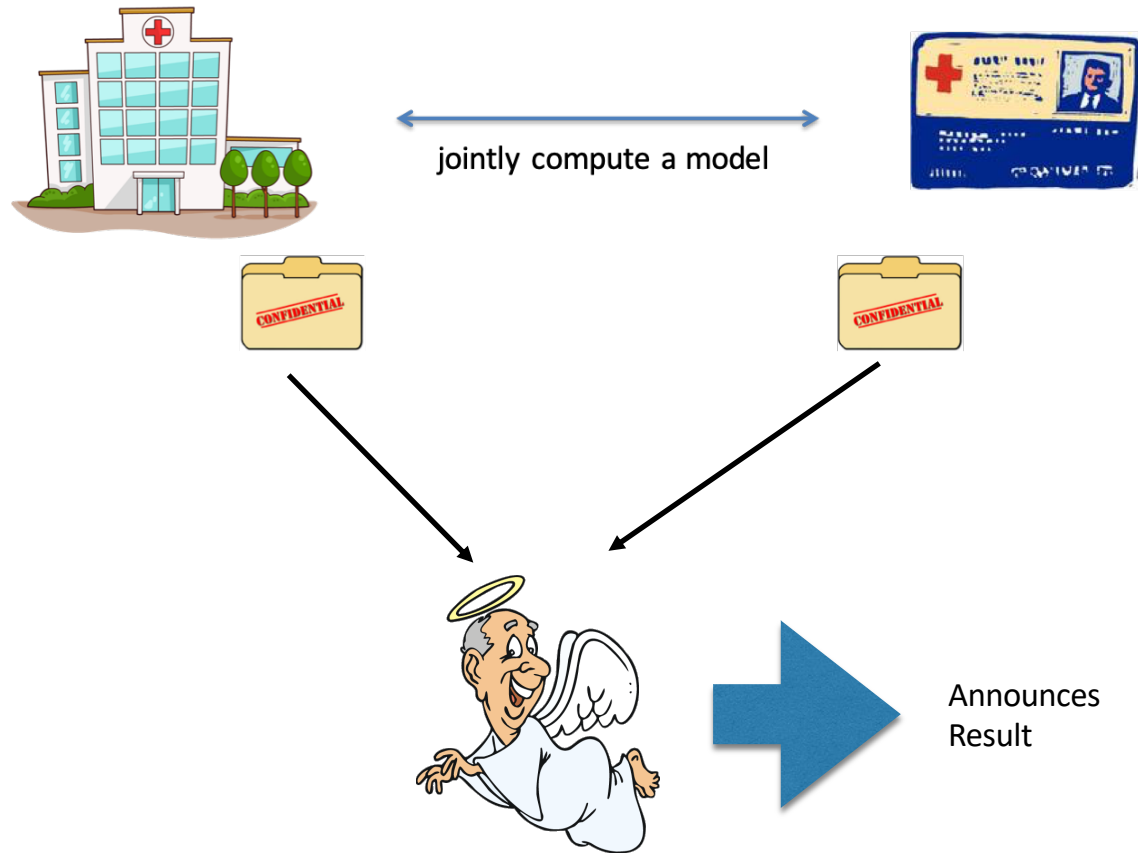Both parties want to guarantee the privacy of their data.

# Crypto Tools

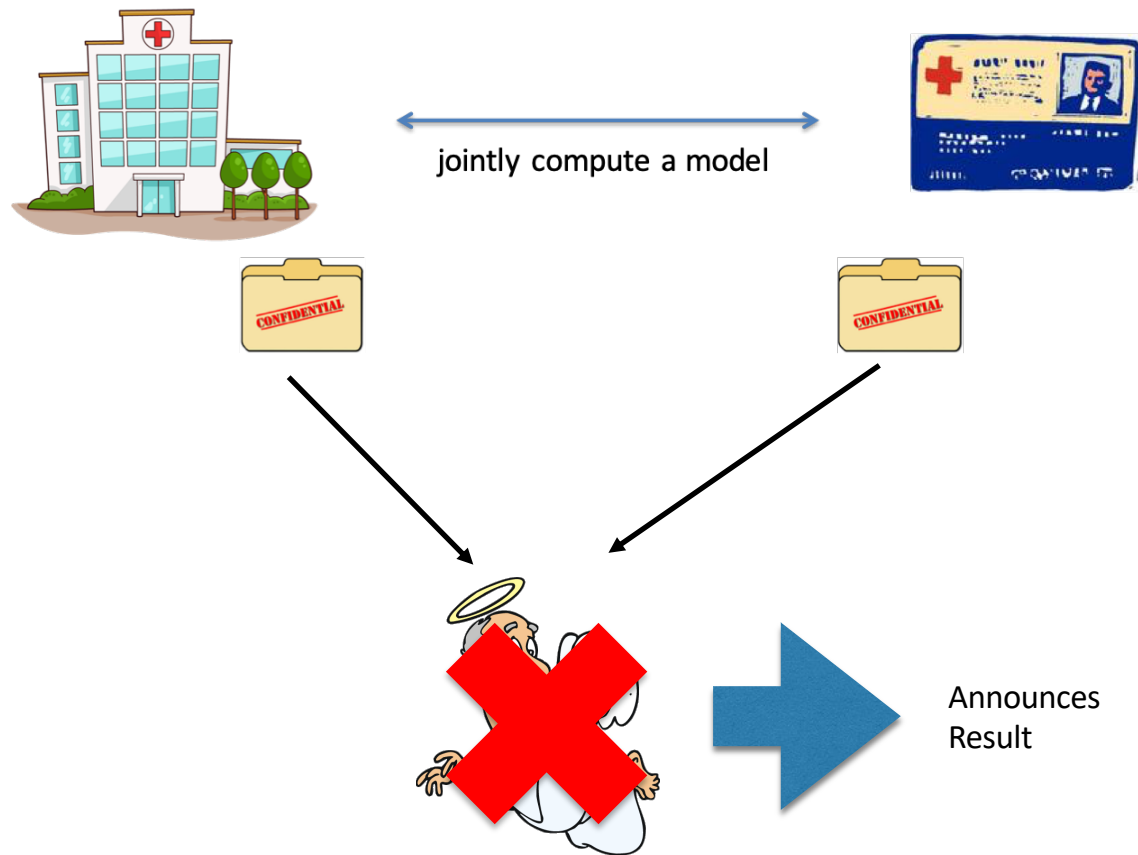Cryptographic tools can be used to obtain privacy guarantees:

◇ Fully homomorphic encryption (FHE)

◇ Differential privacy (DP) mechanisms

◇ Secure multi-party computation (MPC)

# Trusted Third Party



jointly compute a model

Announces Result
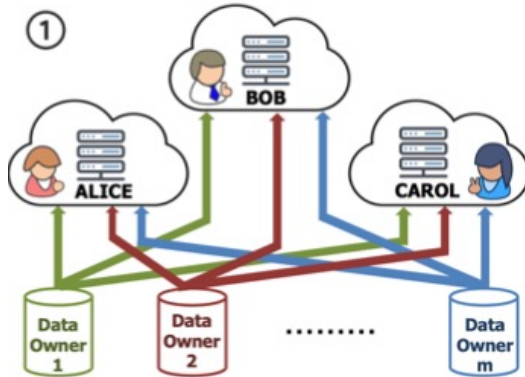
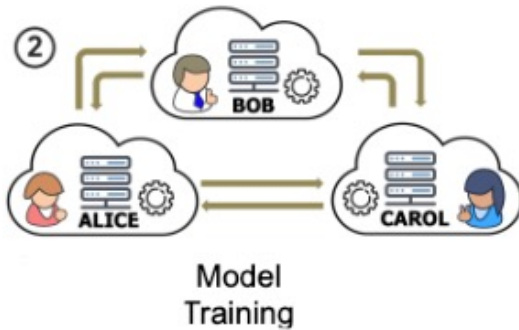# Secure MPC



jointly compute a model

Announces
Result

With MPC it is possible to securely perform the computation without TTP. The parties just learn the desired output, and nothing more.
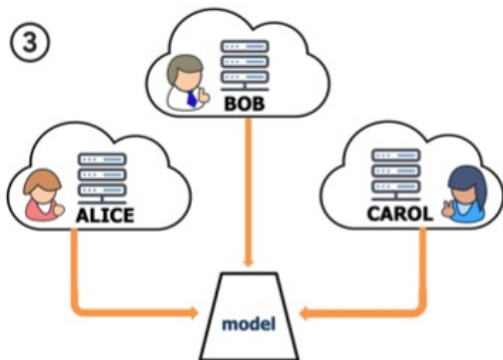
# Computation Using Secret Shares



Data owners convert input data to integers modulo q and secret share it with the computing parties.

Computation is securely performed over the secret shared values.

Secret shares corresponding to the outputs are reconstructed towards the parties that should get the output of the computation.

# Efficient Solutions

In principle, functions can be securely computed using MPC by expressing them in terms of addition and multiplications gates (alternatively, AND and XOR gates).

For efficiency reasons, it is desirable to optimize many of the important operations (e.g., secure comparison, secure linear algebra operations, etc). And operations can be more efficient in either $Z_2$ or $Z_q$ for large q.

Setting: Malicious or semi-honest adversaries? How many parties? How many corruptions? This influences the underlying MPC scheme that should be used.

To get very efficient PPML solutions, it is often the case that the ML techniques need to be made more "**MPC-friendly**".

# Using MPC for PPML

PPML based on MPC is a hot area with many works appearing in top venues in recent years, dealing with topics such as:

- Training and scoring of neural networks

- Linear regression

- Training of tree ensembles (with continuous attributes)

- Privacy-preserving video classification

- …

Also a lot of progress in the underlying MPC schemes in the last few years.

# Scalable Framework for PPML

Privacy-preserving machine learning solutions so far mostly assume that the datasets needed for training are ready for consumption.

In reality much work is done before the training:

- Datasets must be cleaned and pre-processed.

- Missing values need to be addressed.

- Continuous data must be discretized.

- Training features must be selected.

Models must be validated and possibly fine-tuned.

# Privacy-Preserving Feature Selection

First MPC based protocol for privacy-preserving feature selection. It is based on the filter method, which is independent of model training, and can be used in combination with any MPC protocol to rank features.

Efficient feature scoring protocol based on Gini impurity.

We show that secure feature selection with the proposed protocols improves the accuracy of classifiers on a variety of real-world data sets, without leaking information about the feature values or even which features were selected.

*Xiling Li, Rafael Dowsley, Martine De Cock. Privacy-Preserving Feature Selection with Secure Multiparty Computation. ICML 2021*

# Feature Selection

Assume a data set $S$ of $m$ training examples, where each training example consists of an input feature vector $(x_1, \ldots, x_p)$ and a corresponding label $y$ (n possible classes).

Not all $p$ features may be equally beneficial for training ML models. Using a well-chosen subset of features can lead to more accurate models, as well as efficiency gains during model training.

In the filter approach to feature selection, all features are first assigned a score that is indicative of their predictive ability. Then only the best scoring features are retained.

Well-known techniques to score features in terms of their informativeness include mutual information (MI), **Gini impurity (GI)**, and Pearson's correlation coefficient (PCC).

# Gini Impurity

The computation of a GI score for continuous valued features traditionally requires sorting of the feature values to determine candidate split points in the feature value range.

As sorting is an expensive operation to perform in a privacy-preserving way, we instead propose a "mean-split Gini score" (MS-GINI) that avoids the need for sorting by selecting the mean of the feature values as the split point.

Feature selection with MS-GINI leads to accuracy improvements that are on par with those obtained with GI, PCC, and MI in the data sets used in our experiments.

# Adversarial Model

Our protocols for privacy-preserving feature selection are sufficiently generic to be used in dishonest-majority as well as honest majority settings, and with passive or active adversaries, just by changing the underlying MPC scheme.

Some of the most efficient MPC schemes are for 3 parties and at most one corruption (e.g., Araki et al., CCS 2016). We evaluate the runtime of our protocols in this honest-majority 3PC setting, which is growing in popularity in the PPML literature.

In the case of malicious adversaries we show how even better runtimes can be obtained with a recently proposed MPC scheme for 4PC with one corruption of Dalskov et al. (Usenix 2021).

# Passive 3PC with Replicated Secret Sharing

To create a replicated secret sharing $[x]$ of a value $x$, pick uniformly random $x_1,x_2,x_3$ subject to the constraint that $x_1+x_2+x_3=x$ mod q.

$x_1 , x_2$

$x_2 , x_3$

$x_3 , x_1$

# Passive 3PC with Replicated Secret Sharing

Any single party cannot learn any information about $x$ given its shares.

$$x_1 + x_2 + x_3 = x \bmod q$$

$x_1, x_2$

$x_2, x_3$

$x_3, x_1$

# Passive 3PC with Replicated Secret Sharing

Addition of secret shared values are <span style="color:red">easily</span> performed locally by the parties.

$$x_1+x_2+x_3=x \bmod q$$

$$y_1+y_2+y_3=y \bmod q$$

$x_1, x_2, y_1, y_2$

$x_2, x_3, y_2, y_3$

$x_3, x_1, y_3, y_1$

# Passive 3PC with Replicated Secret Sharing

Addition of secret shared values are easily performed locally by the parties.

Replicated secret sharing of $x+y$

$x_1+y_1$, $x_2+y_2$

$x_2+y_2$, $x_3+y_3$

$x_3+y_3$, $x_1+y_1$

It is also easy to add a constant or multiply by a constant.

# Passive 3PC with Replicated Secret Sharing

Main advantage of replicated secret sharing: a very efficient procedure for doing one level of multiplication of secret shared values.

$$x_1 + x_2 + x_3 = x \bmod q$$

$$y_1 + y_2 + y_3 = y \bmod q$$

$x_1, x_2, y_1, y_2$

$x_2, x_3, y_2, y_3$

$x_3, x_1, y_3, y_1$

# Passive 3PC with Replicated Secret Sharing

Main advantage of replicated secret sharing: a very efficient procedure for doing one level of multiplication of secret shared values.
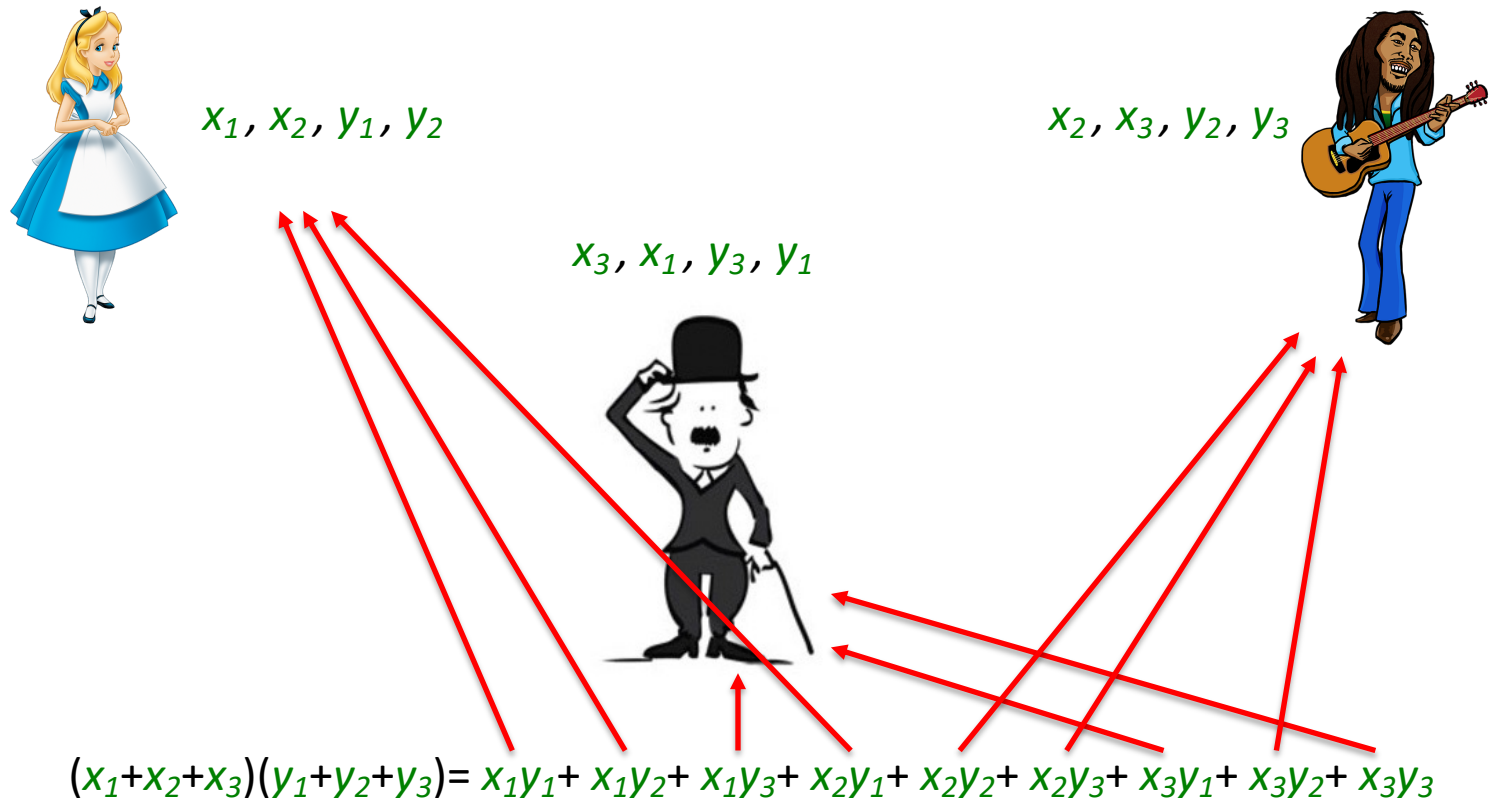
$x_1, x_2, y_1, y_2$

$x_2, x_3, y_2, y_3$

$x_3, x_1, y_3, y_1$

$$(x_1+x_2+x_3)(y_1+y_2+y_3)= x_1y_1+ x_1y_2+ x_1y_3+ x_2y_1+ x_2y_2+ x_2y_3+ x_3y_1+ x_3y_2+ x_3y_3$$

# Passive 3PC with Replicated Secret Sharing

Main advantage of replicated secret sharing: a very efficient procedure for doing one level of multiplication of secret shared values.

These terms can be locally computed by Alice.

$x_1, x_2, y_1, y_2$

$x_2, x_3, y_2, y_3$

$x_3, x_1, y_3, y_1$

$(x_1+x_2+x_3)(y_1+y_2+y_3) = x_1y_1 + x_1y_2 + x_1y_3 + x_2y_1 + x_2y_2 + x_2y_3 + x_3y_1 + x_3y_2 + x_3y_3$

# Passive 3PC with Replicated Secret Sharing

Main advantage of replicated secret sharing: a very efficient procedure for doing one level of multiplication of secret shared values.

$x_1, x_2, y_1, y_2$

$x_2, x_3, y_2, y_3$

$x_3, x_1, y_3, y_1$

$$(x_1+x_2+x_3)(y_1+y_2+y_3)= x_1y_1+ x_1y_2+ x_1y_3+ x_2y_1+ x_2y_2+ x_2y_3+ x_3y_1+ x_3y_2+ x_3y_3$$

# Passive 3PC with Replicated Secret Sharing

Main advantage of replicated secret sharing: a very efficient procedure for doing one level of multiplication of secret shared values.

$$z_1+z_2+z_3=xy \bmod q$$

$$z_1 = x_1y_1 + x_1y_2 + x_2y_1 \qquad\qquad z_2 = x_2y_2 + x_2y_3 + x_3y_2$$

$$z_3 = x_1y_3 + x_3y_1 + x_3y_3$$

The parties only obtain an additively secret sharing of $xy$.

# Passive 3PC with Replicated Secret Sharing

To convert back to the replicated secret sharing format, the parties add a secret sharing of zero (can be generated using PRFs) to the secret sharing of z to obtain $u_1$, $u_2$, $u_3$ and then exchange the appropriate shares.

$u_1 + u_2 + u_3 = xy$ mod q

$u_1$

$u_2$

$u_3$

# Passive 3PC with Replicated Secret Sharing

To convert back to the replicated secret sharing format, the parties add a secret sharing of zero (can be generated using PRFs) to the secret sharing of z to obtain $u_1$, $u_2$, $u_3$ and then exchange the appropriate shares.

# Passive 3PC with Replicated Secret Sharing

To convert back to the replicated secret sharing format, the parties add a secret sharing of zero (can be generated using PRFs) to the secret sharing of z to obtain $u_1$, $u_2$, $u_3$ and then exchange the appropriate shares.

$u_1 + u_2 + u_3 = xy$ mod q

$u_1$, $u_2$

$u_2$, $u_3$

$u_3$, $u_1$

# Active 3PC with Replicated Secret Sharing

In the active 3PC setting, information-theoretic message authentication codes (MACs) are used to prevent the parties from deviating from the protocol specification. The MAC is similar to the one from SPDZ.

In addition to computations over secret shares of the data, the parties also perform computations required for MACs.

# Active 4PC with 1 Corruption

$x_1 + x_2 + x_3 + x_4 = x \bmod q$



$x_4, x_1, x_2$

$x_1, x_2, x_3$

$x_2, x_3, x_4$

$x_3, x_4, x_1$

In active 4PC setting, we use Fantastic Four (Dalskov et al., Usenix 2021).

# Active 4PC with 1 Corruption



**Joint message passing protocol** allows 2 parties holding the same input $u$ to send $u$ to a third party.

# Active 4PC with 1 Corruption



*u*

*u*

In case something goes wrong, one or two parties are identified (one of which is malicious). One of the identified parties is excluded, the shares converted to replicated secret sharing for 3PC, and the execution continues with 3PC.

# Active 4PC with 1 Corruption

$y_1 + y_2 + y_3 = x \bmod q$



$y_2, y_3$

$y_1, y_2$

$y_3, y_1$

If problems persist, the malicious party is identified with certainty, and the shares converted to additively secret sharing for 2PC.

# Active 4PC with 1 Corruption

$z_1 + z_2 = x$ mod q

$z_2$

$z_1$

The first excluded party can now help generating multiplication triples for the 2PC.

If problems persist, the malicious party is identified with certainty, and the shares converted to additively secret sharing for 2PC.

# Building Blocks

- Secure comparison protocol: Given [*x*] and [*y*], returns [*1*] if *x* < *y* and [*0*] otherwise.

- Secure argmin protocol: Given a secret shared vector, returns a secret sharing of the index at which the vector has the minimum value.

- Secure equality protocol: Given [*x*] and [*y*], returns [*1*] if *x* = *y* and [*0*] otherwise.

- Secure division protocol: Given [*x*] and [*y*], returns [*z*] where *z* = *x*/*y*.

# Filter-Based Feature Selection

**Protocol 1** Protocol $\pi_{\mathsf{FILTER-FS}}$ for Secure Filter based Feature Selection

**Input:** A secret shared $m \times p$ data matrix $[\![D]\!]_q$, a secret shared $p$-length score vector $[\![G]\!]_q$, the number $k < p$ of features to be selected, and a constant $t$ that is bigger than the highest possible score in $[\![G]\!]_q$

**Output:** a secret shared $m \times k$ matrix $[\![D']\!]_q$

1: **for** $i = 1$ **to** $k$ **do**
2:     $[\![I[i]]\!]_q \leftarrow \pi_{\mathsf{ARGMIN}}([\![G]\!]_q)$
3:     **for** $j \leftarrow 1$ **to** $p$ **do**
4:         $[\![flag_k]\!]_q \leftarrow \pi_{\mathsf{EQ}}([\![I[i]]\!]_q, j)$
5:         $[\![T[j][i]]\!]_q \leftarrow [\![flag_k]\!]_q$
6:         $[\![G[j]]\!]_q \leftarrow [\![G[j]]\!]_q + \pi_{\mathsf{DM}}([\![flag_k]\!]_q, t - [\![G[j]]\!]_q)$
7:     **end for**
8: **end for**
9: $[\![D']\!]_q \leftarrow \pi_{\mathsf{DMM}}([\![D]\!]_q, [\![T]\!]_q)$
10: **return** $[\![D']\!]_q$

Inputs:
- Data matrix $D$ of size $m$ x $p$
- A vector $G$ with the scores of the $p$ features
- Upper bound $t$ on the scores

Output:
- Data matrix $D'$ of size $m$ x $k$ containing the columns of $D$ corresponding to the $k$ lowest scores in $G$.

# Feature Scoring Using MS-GINI

Split the set of values of feature $F_j$ based on its mean value as a threshold $\theta$.

Let $S_{\leq\theta}$ denote the set of instances that have $x_j \leq \theta$, $S_{>\theta}$ denote the set of instances that have $x_j > \theta$ and $L_c$ the set of examples from $S$ that have class label $c$. Define:

$$p_c^{\leq\theta} = \frac{|S_{\leq\theta} \cap L_c|}{|S_{\leq\theta}|}; \quad p_c^{>\theta} = \frac{|S_{>\theta} \cap L_c|}{|S_{>\theta}|}$$

Then we get the Gini impurities:

$$G(S_{\leq\theta}) = 1 - \sum_{c=1}^{n}(p_c^{\leq\theta})^2; \quad G(S_{>\theta}) = 1 - \sum_{c=1}^{n}(p_c^{>\theta})^2$$

and the Gini score:

$$G(F_j) = \frac{1}{m} \cdot (|S_{\leq\theta}| \cdot G(S_{\leq\theta}) + |S_{>\theta}| \cdot G(S_{>\theta}))$$

# Feature Scoring Using MS-GINI

---

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

---

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q$, $[\![b]\!]_q$, $[\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:     $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:     $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:     **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:         $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:         $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:         $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:    **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}} ( \pi_{\mathsf{DP}} ([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}} (1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}} ( \pi_{\mathsf{DP}} ([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}} (1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

---

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q, [\![b]\!]_q, [\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:    $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:    $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:    **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:       $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:       $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:       $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:   **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}} ( \pi_{\mathsf{DP}} ([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}} (1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}} ( \pi_{\mathsf{DP}} ([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}} (1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

Inputs:
- Secret shares of a feature column *F*
- Secret shares of one-hot-encoded version of the label vector *L*

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

---

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$ ⬅

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + \dots + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q$, $[\![b]\!]_q$, $[\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:      $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:      $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:      **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:          $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:          $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:          $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:      **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + \dots + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + \dots + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}}(1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}}(1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

---

Inputs:
- Secret shares of a feature column *F*
- Secret shares of one-hot-encoded version of the label vector *L*

Output:
- Secret share of the score of *F*

# Feature Scoring Using MS-GINI

---

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

---

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$   $\longleftarrow$
2: Initialize $[\![a]\!]_q, [\![b]\!]_q, [\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:    $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:    $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:    **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:       $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:       $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:       $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:    **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}}(1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}}(1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

---

Compute the split point as the mean of the feature values in the column.

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

**Input:** A secret shared feature column $\llbracket F \rrbracket_q = (\llbracket f_1 \rrbracket_q, \llbracket f_2 \rrbracket_q, \ldots, \llbracket f_m \rrbracket_q)$, a secret shared $m \times (n-1)$ label-class matrix $\llbracket L \rrbracket_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $\llbracket G(F) \rrbracket_q$ of the feature $F$

1: $\llbracket \theta \rrbracket_q \leftarrow (\llbracket f_1 \rrbracket_q + \llbracket f_2 \rrbracket_q + \ldots + \llbracket f_m \rrbracket_q) \cdot \frac{1}{m}$
2: Initialize $\llbracket a \rrbracket_q$, $\llbracket b \rrbracket_q$, $\llbracket A \rrbracket_q$ and $\llbracket B \rrbracket_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:     $\llbracket flag_s \rrbracket_q \leftarrow \pi_{\mathsf{LT}}(\llbracket \theta \rrbracket_q, \llbracket f_i \rrbracket_q)$
5:     $\llbracket b \rrbracket_q \leftarrow \llbracket b \rrbracket_q + \llbracket flag_s \rrbracket_q$
6:     **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:         $\llbracket flag_m \rrbracket_q \leftarrow \pi_{\mathsf{DM}}(\llbracket flag_s \rrbracket_q, \llbracket L[i][j] \rrbracket_q)$
8:         $\llbracket B[j] \rrbracket_q \leftarrow \llbracket B[j] \rrbracket_q + \llbracket flag_m \rrbracket_q$
9:         $\llbracket A[j] \rrbracket_q \leftarrow \llbracket A[j] \rrbracket_q + \llbracket L[i][j] \rrbracket_q - \llbracket flag_m \rrbracket_q$
10:    **end for**
11: **end for**
12: $\llbracket a \rrbracket_q \leftarrow m - \llbracket b \rrbracket_q$
13: $\llbracket A[n] \rrbracket_q \leftarrow \llbracket a \rrbracket_q - (\llbracket A[1] \rrbracket_q + \ldots + \llbracket A[n-1] \rrbracket_q)$
14: $\llbracket B[n] \rrbracket_q \leftarrow \llbracket b \rrbracket_q - (\llbracket B[1] \rrbracket_q + \ldots + \llbracket B[n-1] \rrbracket_q)$
15: $\llbracket G(S_{\leq \theta}) \rrbracket_q \leftarrow \llbracket a \rrbracket_q - \pi_{\mathsf{DM}} (\pi_{\mathsf{DP}} (\llbracket A \rrbracket_q, \llbracket A \rrbracket_q), \pi_{\mathsf{DIV}} (1, \llbracket a \rrbracket_q))$
16: $\llbracket G(S_{> \theta}) \rrbracket_q \leftarrow \llbracket b \rrbracket_q - \pi_{\mathsf{DM}} (\pi_{\mathsf{DP}} (\llbracket B \rrbracket_q, \llbracket B \rrbracket_q), \pi_{\mathsf{DIV}} (1, \llbracket b \rrbracket_q))$
17: $\llbracket G(F) \rrbracket_q \leftarrow \llbracket G(S_{\leq \theta}) \rrbracket_q + \llbracket G(S_{> \theta}) \rrbracket_q$
18: **return** $\llbracket G(F) \rrbracket_q$

Lines 2-14 have the goal of getting the following values in the counters $a$, $b$, $A[j]$, $B[j]$:

$$
\begin{aligned}
a &= |S_{\leq \theta}| \\
b &= |S_{> \theta}| \\
A[j] &= |S_{\leq \theta} \cap L_j|, \text{ for } j = 1 \ldots n \\
B[j] &= |S_{> \theta} \cap L_j|, \text{ for } j = 1 \ldots n
\end{aligned}
$$

# Feature Scoring Using MS-GINI

---

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

---

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

Initialize $a, b, A[j], B[j]$ with zeros.

1:  $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2:  Initialize $[\![a]\!]_q, [\![b]\!]_q, [\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.    ←
3:  **for** $i \leftarrow 1$ **to** $m$ **do**
4:      $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:      $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:      **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:          $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:          $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:          $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:     **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}}(1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}}(1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

---

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q$, $[\![b]\!]_q$, $[\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:   $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:   $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:   **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:     $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:     $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:     $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:    **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}}(1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}}(1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

Perform a secure comparison to determine whether the instance belongs to $S_{>\theta}$.

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

---

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q, [\![b]\!]_q, [\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:      $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:      $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$     ⬅
6:      **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:          $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:          $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:          $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:      **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}}(1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}}(1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

---

Perform a secure comparison to determine whether the instance belongs to $S_{>\theta}$.

The value of counter $b$ is updated accordingly, and in an oblivious way.

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

---

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, \dots, [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + \dots + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q$, $[\![b]\!]_q$, $[\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:     $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:     $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:     **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:         $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:         $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:         $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:    **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + \dots + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + \dots + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}} (\pi_{\mathsf{DP}} ([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}} (1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}} (\pi_{\mathsf{DP}} ([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}} (1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

---

Perform a secure comparison to determine whether the instance belongs to $S_{>\theta}$.

The value of counter *b* is updated accordingly, and in an oblivious way.

Counter *a* will be computed in an indirect way later.

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

---

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q$, $[\![b]\!]_q$, $[\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:     $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:     $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:     **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:         $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$   ⬅
8:         $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$   ⬅
9:         $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:    **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}} (\pi_{\mathsf{DP}} ([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}} (1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}} (\pi_{\mathsf{DP}} ([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}} (1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

---

Perform a secure comparison to determine whether the instance belongs to $S_{>\theta}$.

The value of counter *b* is updated accordingly, and in an oblivious way.

Check whether the instance belongs to $S_{>\theta} \cap L_j$ and obliviously update *B[j]* accordingly.

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q$, $[\![b]\!]_q$, $[\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4: $\quad [\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5: $\quad [\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6: $\quad$ **for** $j \leftarrow 1$ **to** $n-1$ **do**
7: $\quad\quad [\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$ $\quad\longleftarrow$
8: $\quad\quad [\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9: $\quad\quad [\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$ $\quad\longleftarrow$
10: $\quad$ **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}}(1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}}(1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

Perform a secure comparison to determine whether the instance belongs to $S_{>\theta}$.

The value of counter $b$ is updated accordingly, and in an oblivious way.

Check whether the instance belongs to $S_{>\theta} \cap L_j$ and obliviously update $B[j]$ accordingly.

Check whether the instance belongs to $S_{\leq\theta} \cap L_j$ and obliviously update $A[j]$ accordingly. $flag_m$ is reused to save on secure multiplications.

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

---

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q$, $[\![b]\!]_q$, $[\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:     $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:     $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:     **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:         $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:         $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:         $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:    **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$   ⬅
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$   ⬅
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}} ( \pi_{\mathsf{DP}} ([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}} (1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}} ( \pi_{\mathsf{DP}} ([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}} (1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

---

Optimization to compute *A[n], B[n]* without needing a n$^{th}$ iteration of the inner **for** loop.

These operations are done locally by each party.

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\text{MS-GINI}}$ for Secure MS-GINI Score of a Feature

---

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q$, $[\![b]\!]_q$, $[\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:    $[\![flag_s]\!]_q \leftarrow \pi_{\text{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:    $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:    **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:       $[\![flag_m]\!]_q \leftarrow \pi_{\text{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:       $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:       $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:   **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\text{DM}}(\pi_{\text{DP}}([\![A]\!]_q, [\![A]\!]_q), \pi_{\text{DIV}}(1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\text{DM}}(\pi_{\text{DP}}([\![B]\!]_q, [\![B]\!]_q), \pi_{\text{DIV}}(1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

---

Using the counters

$$a = |S_{\leq\theta}|$$
$$b = |S_{>\theta}|$$
$$A[j] = |S_{\leq\theta} \cap L_j|, \text{ for } j = 1 \ldots n$$
$$B[j] = |S_{>\theta} \cap L_j|, \text{ for } j = 1 \ldots n$$

the Gini score can be computed as:

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q$, $[\![b]\!]_q$, $[\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:     $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:     $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:     **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:         $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:         $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:         $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:     **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}}(1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}}(1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

Using the counters

$$
\begin{aligned}
a &= |S_{\leq\theta}| \\
b &= |S_{>\theta}| \\
A[j] &= |S_{\leq\theta} \cap L_j|, \text{ for } j = 1 \ldots n \\
B[j] &= |S_{>\theta} \cap L_j|, \text{ for } j = 1 \ldots n
\end{aligned}
$$

the Gini score can be computed as:

$$
\begin{aligned}
G(F) &= \frac{1}{m} \cdot \left[ a \cdot \left( 1 - \sum_{j=1}^{n} \left( \frac{A[j]}{a} \right)^2 \right) + b \cdot \left( 1 - \sum_{j=1}^{n} \left( \frac{B[j]}{b} \right)^2 \right) \right] \\
&= \frac{1}{m} \cdot \left[ \left( a - \frac{1}{a} \cdot A \bullet A \right) + \left( b - \frac{1}{b} \cdot B \bullet B \right) \right]
\end{aligned}
$$

where the "big dot" denotes dot product.

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

**Input:** A secret shared feature column $[\![F]\!]_q = ([\![f_1]\!]_q, [\![f_2]\!]_q, ..., [\![f_m]\!]_q)$, a secret shared $m \times (n-1)$ label-class matrix $[\![L]\!]_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $[\![G(F)]\!]_q$ of the feature $F$

1: $[\![\theta]\!]_q \leftarrow ([\![f_1]\!]_q + [\![f_2]\!]_q + ... + [\![f_m]\!]_q) \cdot \frac{1}{m}$
2: Initialize $[\![a]\!]_q$, $[\![b]\!]_q$, $[\![A]\!]_q$ and $[\![B]\!]_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:    $[\![flag_s]\!]_q \leftarrow \pi_{\mathsf{LT}}([\![\theta]\!]_q, [\![f_i]\!]_q)$
5:    $[\![b]\!]_q \leftarrow [\![b]\!]_q + [\![flag_s]\!]_q$
6:    **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:      $[\![flag_m]\!]_q \leftarrow \pi_{\mathsf{DM}}([\![flag_s]\!]_q, [\![L[i][j]]\!]_q)$
8:      $[\![B[j]]\!]_q \leftarrow [\![B[j]]\!]_q + [\![flag_m]\!]_q$
9:      $[\![A[j]]\!]_q \leftarrow [\![A[j]]\!]_q + [\![L[i][j]]\!]_q - [\![flag_m]\!]_q$
10:   **end for**
11: **end for**
12: $[\![a]\!]_q \leftarrow m - [\![b]\!]_q$
13: $[\![A[n]]\!]_q \leftarrow [\![a]\!]_q - ([\![A[1]]\!]_q + ... + [\![A[n-1]]\!]_q)$
14: $[\![B[n]]\!]_q \leftarrow [\![b]\!]_q - ([\![B[1]]\!]_q + ... + [\![B[n-1]]\!]_q)$
15: $[\![G(S_{\leq\theta})]\!]_q \leftarrow [\![a]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![A]\!]_q, [\![A]\!]_q), \pi_{\mathsf{DIV}}(1, [\![a]\!]_q))$
16: $[\![G(S_{>\theta})]\!]_q \leftarrow [\![b]\!]_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}([\![B]\!]_q, [\![B]\!]_q), \pi_{\mathsf{DIV}}(1, [\![b]\!]_q))$
17: $[\![G(F)]\!]_q \leftarrow [\![G(S_{\leq\theta})]\!]_q + [\![G(S_{>\theta})]\!]_q$
18: **return** $[\![G(F)]\!]_q$

Using the counters

$$
\begin{aligned}
a &= |S_{\leq\theta}| \\
b &= |S_{>\theta}| \\
A[j] &= |S_{\leq\theta} \cap L_j|, \text{ for } j = 1 \ldots n \\
B[j] &= |S_{>\theta} \cap L_j|, \text{ for } j = 1 \ldots n
\end{aligned}
$$

the Gini score can be computed as:

$$
\begin{aligned}
G(F) &= \frac{1}{m} \cdot \left[ a \cdot \left( 1 - \sum_{j=1}^{n} \left( \frac{A[j]}{a} \right)^2 \right) + b \cdot \left( 1 - \sum_{j=1}^{n} \left( \frac{B[j]}{b} \right)^2 \right) \right] \\
&= \frac{1}{m} \cdot \left[ \left( a - \frac{1}{a} \cdot A \bullet A \right) + \left( b - \frac{1}{b} \cdot B \bullet B \right) \right]
\end{aligned}
$$

where the "big dot" denotes dot product.

# Feature Scoring Using MS-GINI

**Protocol 2** Protocol $\pi_{\mathsf{MS-GINI}}$ for Secure MS-GINI Score of a Feature

---

**Input:** A secret shared feature column $\llbracket F \rrbracket_q = (\llbracket f_1 \rrbracket_q, \llbracket f_2 \rrbracket_q, ..., \llbracket f_m \rrbracket_q)$, a secret shared $m \times (n-1)$ label-class matrix $\llbracket L \rrbracket_q$, where $m$ is the number of instances and $n$ is the number of classes.

**Output:** MS-GINI score $\llbracket G(F) \rrbracket_q$ of the feature $F$

1: $\llbracket \theta \rrbracket_q \leftarrow (\llbracket f_1 \rrbracket_q + \llbracket f_2 \rrbracket_q + ... + \llbracket f_m \rrbracket_q) \cdot \frac{1}{m}$
2: Initialize $\llbracket a \rrbracket_q$, $\llbracket b \rrbracket_q$, $\llbracket A \rrbracket_q$ and $\llbracket B \rrbracket_q$ with zeros.
3: **for** $i \leftarrow 1$ **to** $m$ **do**
4:     $\llbracket flag_s \rrbracket_q \leftarrow \pi_{\mathsf{LT}}(\llbracket \theta \rrbracket_q, \llbracket f_i \rrbracket_q)$
5:     $\llbracket b \rrbracket_q \leftarrow \llbracket b \rrbracket_q + \llbracket flag_s \rrbracket_q$
6:     **for** $j \leftarrow 1$ **to** $n-1$ **do**
7:         $\llbracket flag_m \rrbracket_q \leftarrow \pi_{\mathsf{DM}}(\llbracket flag_s \rrbracket_q, \llbracket L[i][j] \rrbracket_q)$
8:         $\llbracket B[j] \rrbracket_q \leftarrow \llbracket B[j] \rrbracket_q + \llbracket flag_m \rrbracket_q$
9:         $\llbracket A[j] \rrbracket_q \leftarrow \llbracket A[j] \rrbracket_q + \llbracket L[i][j] \rrbracket_q - \llbracket flag_m \rrbracket_q$
10:    **end for**
11: **end for**
12: $\llbracket a \rrbracket_q \leftarrow m - \llbracket b \rrbracket_q$
13: $\llbracket A[n] \rrbracket_q \leftarrow \llbracket a \rrbracket_q - (\llbracket A[1] \rrbracket_q + ... + \llbracket A[n-1] \rrbracket_q)$
14: $\llbracket B[n] \rrbracket_q \leftarrow \llbracket b \rrbracket_q - (\llbracket B[1] \rrbracket_q + ... + \llbracket B[n-1] \rrbracket_q)$
15: $\llbracket G(S_{\leq \theta}) \rrbracket_q \leftarrow \llbracket a \rrbracket_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}(\llbracket A \rrbracket_q, \llbracket A \rrbracket_q), \pi_{\mathsf{DIV}}(1, \llbracket a \rrbracket_q))$ ⬅
16: $\llbracket G(S_{> \theta}) \rrbracket_q \leftarrow \llbracket b \rrbracket_q - \pi_{\mathsf{DM}}(\pi_{\mathsf{DP}}(\llbracket B \rrbracket_q, \llbracket B \rrbracket_q), \pi_{\mathsf{DIV}}(1, \llbracket b \rrbracket_q))$ ⬅
17: $\llbracket G(F) \rrbracket_q \leftarrow \llbracket G(S_{\leq \theta}) \rrbracket_q + \llbracket G(S_{> \theta}) \rrbracket_q$ ⬅
18: **return** $\llbracket G(F) \rrbracket_q$

---

Using the counters

$$
\begin{aligned}
a &= |S_{\leq \theta}| \\
b &= |S_{> \theta}| \\
A[j] &= |S_{\leq \theta} \cap L_j|, \text{ for } j = 1 \ldots n \\
B[j] &= |S_{> \theta} \cap L_j|, \text{ for } j = 1 \ldots n
\end{aligned}
$$

the Gini score can be computed as:

$$
\begin{aligned}
G(F) &= \frac{1}{m} \cdot \left[ a \cdot \left( 1 - \sum_{j=1}^{n} \left( \frac{A[j]}{a} \right)^2 \right) + b \cdot \left( 1 - \sum_{j=1}^{n} \left( \frac{B[j]}{b} \right)^2 \right) \right] \\
&= \frac{1}{m} \cdot \left[ \left( a - \frac{1}{a} \cdot A \bullet A \right) + \left( b - \frac{1}{b} \cdot B \bullet B \right) \right]
\end{aligned}
$$

where the "big dot" denotes dot product.

Factor 1/m is omitted as it has no effect on the relative ordering of the scores of the individual features.

# Accuracy

Data sets:
- Cognitive Load Detection (COG) (Gjoreski et al., 2020)
- Lee Silverman Voice Treatment (LSVT) (Tsanas et al., 2014)
- Speed Dating (SPEED) (Fisman et al., 2006)

*Table 1.* Accuracies of logistic regression models with different feature selection criteria

| Data set | data set details | | | | logistic regression accuracy results | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | $m$ | $p$ | $k$ | #folds | RAW | **MS-GINI (Ours)** | GI | PCC | MI |
| COG | 632 | 120 | 12 | 6 | 50.90% | 52.50% | 52.70% | 48.57% | 51.59% |
| LSVT | 126 | 310 | 103 | 10 | 80.09% | 86.15% | 82.74% | 78.89% | 85.38% |
| SPEED | 8,378 | 122 | 67 | 10 | 95.24% | 97.26% | 95.56% | 95.89% | 95.83% |

As the results show, feature selection based on MS-GINI is on par with the other methods, and substantially improves the accuracy compared to model training on the RAW data sets.

# Runtime

Table 2. Runtime results for privacy-preserving feature selection

| Data set | data set details | | | | **MS-GINI runtime (Ours)** | | | GI runtime (Sorting-based approach) | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | $m$ | $p$ | $k$ | #folds | passive 3PC | active 3PC | active 4PC | passive 3PC | active 3PC | active 4PC |
| COG | 632 | 120 | 12 | 6 | 50 sec | 163 sec | 79 sec | 565 sec | 2,884 sec | 702 sec |
| LSVT | 126 | 310 | 103 | 10 | 60 sec | 254 sec | 89 sec | 368 sec | 1,269 sec | 442 sec |
| SPEED | 8,378 | 122 | 67 | 10 | 949 sec | 3,634 sec | 1,435 sec | 12,241 sec | 97,871 sec | 14,114 sec |

Implemented in MP-SPDZ (Keller, CCS 2020) using ring with $q=2^{64}$.

All benchmark tests were completed on 3 or 4 co-located F32s V2 Azure virtual machines. Each VM contains 32 cores, 64 GiB of memory, and up to a 14 Gbps network bandwidth between each virtual machine.

The reported runtimes are an average across the folds (and covers communication time).

To empirically verify the runtime improvements that can be obtained with our MS-GINI criterion compared to traditional GI, we replaced Protocol 2 by an MPC protocol for computing GI proposed by (Abspoel et al., Usenix 2021).
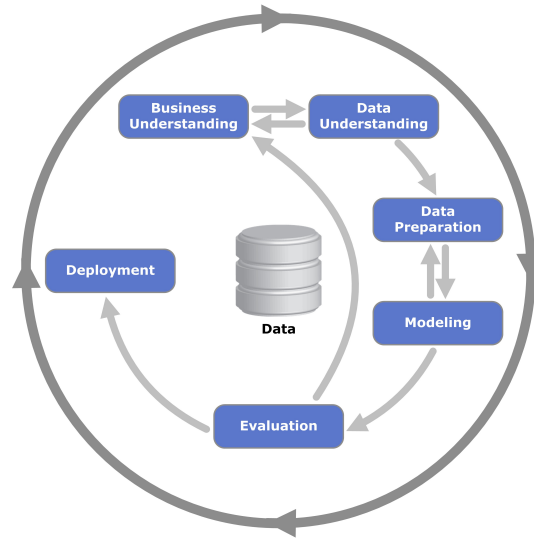
# Future Directions

Development of privacy-preserving protocols for other feature scoring techniques.

MPC protocols for other tasks related to the data preprocessing phase still need to be developed, e.g., privacy-preserving solutions to deal with outliers and missing values.

# Scalable Framework for PPML

Goal: Develop a practical end-to-end framework for privacy-preserving machine learning that covers all the steps necessary for realistic large-scale applications.



Few works tried to close this cycle, and this state of affairs restricts the use of proposed solutions in practical applications.

# MPC+DP

Exploring benefits from combining techniques from MPC and differential privacy is also an interesting direction.