

Enhancing Privacy in Remote Data Classification

A. Piva, C. Orlandi ^{*}, M. Caini, T. Bianchi, and M. Barni

Abstract Neural networks are a fundamental tool in data classification. A protocol whereby a user may ask a service provider to run a neural network on an input provided in encrypted format is proposed here, in such a way that the neural network owner does not get any knowledge about the processed data. At the same time, the knowledge embedded within the network itself is protected. With respect to previous works in this field, the interaction between the user and the neural network owner is kept to a minimum without resorting to general secure multi-party computation protocols.

1 Introduction

Recent advances in signal and information processing together with the possibility of exchanging and transmitting data through flexible and ubiquitous transmission media such as internet and wireless networks, have opened the way towards a new kind of services whereby a provider sells its ability to process and interpret data remotely, e.g. through an internet web service. Examples in this sense include access to remote databases, processing of personal data, processing of multimedia documents, interpretation of medical data for remote diagnosis. In this last scenario, a patient may need a diagnosis from a remote medical institute that has the knowledge needed to perform the diagnosis. Health-related data are of course sensitive, and the patient may do not want to let the institute to know the data he owns; on the other hand, the medical institute is interested in protecting his expertise.

In 2000 two papers [19, 1] proposed the notion of privacy preserving data mining, meaning the possibility to perform data analysis on a distributed database, under some privacy constraints. After the publication of these papers, security constraints were added to several machine learning techniques: decision trees [19], neural networks [6], support vector machines [18], naive bayes classifiers [16], belief networks [24], clustering [14]. In all these works, we can identify two major scenarios: in the first one Alice and Bob share a dataset and want to extract knowledge from it without revealing their own data (we define this scenario as privacy preserving data mining). In the other scenario, which is the one considered in this paper, Alice owns her private data x , while Bob owns an evaluation function C , where in most cases C is a classifier (we define it a remote data classification). Alice is interested in having her data processed by Bob, but she does not want that Bob learns either her input or the output of the computation. At the

Alessandro Piva, Michele Caini, Tiziano Bianchi
Department of Electronics and Telecommunications, University of Florence e-mail: {piva, caini, bianchi}@lci.det.unifi.it

Claudio Orlandi
Department of Computer Science, University of Aarhus, e-mail: orlandi@daimi.au.dk

Mauro Barni
Department of Information Engineering, University of Siena e-mail: barni@dii.unisi.it

^{*} Work done while working at University of Florence.

same time, Bob does not want to reveal the exact form of C , representing his knowledge, since, for instance, he sells a classification service through the web (as in the remote medical diagnosis example).

A fundamental tool in data classification is represented by neural networks (NNs), because of their approximation and generalization capabilities. For this reason, it can be of interest to design a protocol whereby a user may ask a service provider to run a neural network on an input provided in encrypted format. Previous works on privacy preserving NN computing are limited to the systems presented in [3, 6]. However, such studies resort extensively to highly inefficient general secure multi-party computation (SMC) [25] for the computation of the non-linear activation functions implemented in the neurons.

This is not the case with our new protocol which does not resort to general SMC for the evaluation of the activation functions. In a nutshell, the protocol has been designed to ensure that the data provided by the user (say Alice), representing the input of the neural network, are completely protected and, at the same time, to not disclose Bob's classifier (the NN). The proposed protocol relies on homomorphic encryption that allows to perform directly in the encrypted domain all the linear computations. For the non linear functions that can not be handled by means of homomorphic encryption, a limited amount of interaction between the NN owner and the user is introduced to delegate the user (say Alice) to perform some of the computations. Comparing our work with previous ones, another advantage can be highlighted: the proposed protocol can handle every kind of feedforward NN (not only simple layered networks), without disclosing neither the number of neurons nor the way they are connected.

This protocol can be seen as an improvement of the one proposed in [3] from both privacy and efficiency side. In particular, in [3] the authors resorted to a general technique to protect the output value of the neurons, while here we introduce ad-hoc techniques (Section 4.3) for the same purpose. Another difference is in the protection of the network topology, that has been improved and now can cope with more general kind of networks. Finally, in this paper we describe also a practical implementation of our system, proving the feasibility of our technique, for a real set-up. As the best of our knowledge, this is the first practical implementation of a protocol for the private evaluation of a neural network based classifier.

The rest of this paper is organized as follows. In Section 2, a brief overview on Neural Networks that can be used with our protocol is given. In Section 3 our scenario will be described, focusing on the privacy constraints and reviewing the properties to be achieved by our protocol. The way the protocol works is described in Section 4. Section 5 is devoted to the experimental results obtained developing a distributed application that runs the protocol. Some concluding remarks are given in Section 6, while in Appendix how to deal with non integer computation will be discussed.

2 Neural Networks

Neural networks have a great ability to model any given function [17, 13]. Moreover neural networks are provided with good learning algorithms, are robust against noise and generalize well on unseen examples. In this section we will introduce several types of network that can be used with our protocol. The notation is consistent with the one in [5], where a detailed treatment of neural networks is given.

2.0.1 Perceptron

The simplest neural network is the *perceptron* (Figure 1 (a)). It can discriminate between two different classes of instance, and its classification function consists of a linear combination of the input variables, the coefficients of which are the parameters of the model. The discriminant function is of the form $a(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ where \mathbf{x} is the input vector, \mathbf{w} is the vector of weights and w_0 is a threshold². The instance \mathbf{x} is assigned to class c_1 if $a(\mathbf{x}) \geq 0$ and to class c_2 if $a(\mathbf{x}) < 0$. This method can easily be extended to the multiclass case using one discriminant function $a_k(x)$ for each class C_k such that $a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x}$.

² From now on we can forget about w_0 simply appending it at the end of the vector and obtaining $a(\mathbf{x}) = [\mathbf{w} \ w_0]^T [\mathbf{x} \ 1]$

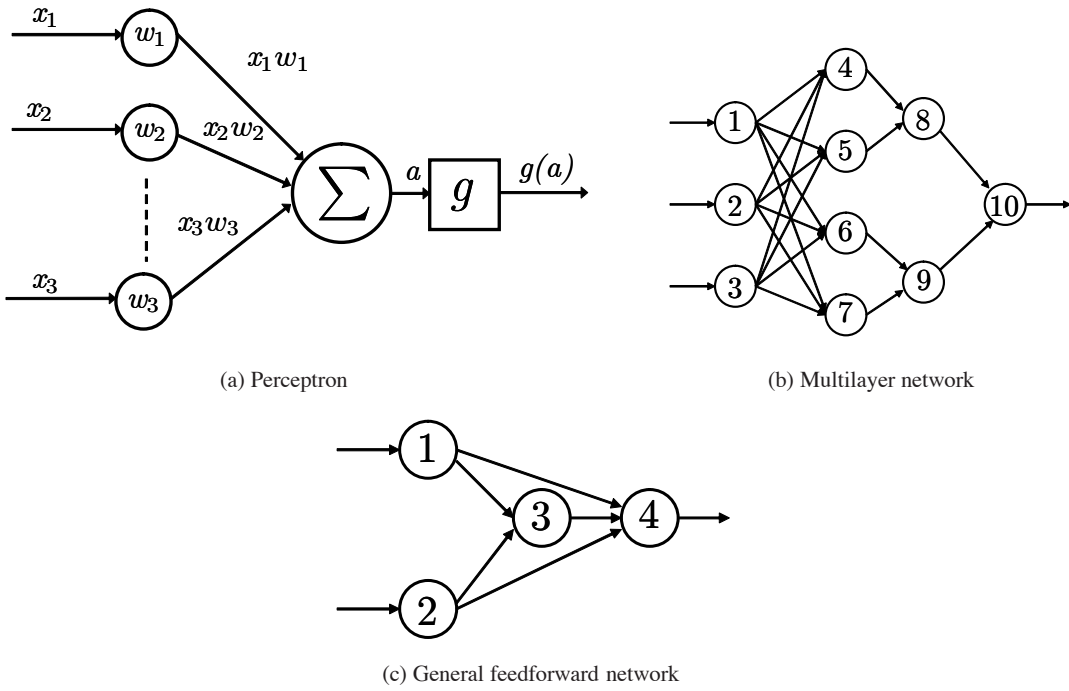


Fig. 1: Three kinds of different networks that can be used with the proposed protocol. Figure (a) shows a single layer network, known as *perceptron*; Figure (b) shows a *multi-layer feedforward network*, while Figure (c) shows a *general feedforward network*. Note that in feed-forward networks it is possible to number neurons such that every neuron gets inputs only from neurons with smaller index.

2.0.2 Feed-Forward Networks

To allow for more general classifications it is possible to consider a network of interconnected neurons. Networks consisting of successive layers of adaptive weights are called *layered networks*: in such a network every unit in one layer is connected to every unit in the next layer, but no other connections are permitted, as shown in Figure 1 (b). The units that are not treated as output units are called *hidden* units.

The output of the j -th hidden unit in the i -th layer is obtained by first forming a weighted linear combination of its l input values to give:

$$a_j^{(i)} = \sum_{k=1}^l w_{jk}^{(i)} x_k$$

Here $w_{jk}^{(i)}$ denotes the weight associated to an edge going from the k -th neuron of the previous layer to the neuron j of i -th layer. The activation of the hidden unit is obtained by transforming the linear sum using an activation function $g(\cdot)$ to obtain $z_j^{(i)} = g(a_j^{(i)})$.

This procedure is iterated until the output layer is reached, obtaining the final output $y_k = g(a_k)$. We will refer to a L -layer network as a network having L layers of adaptive weights, regardless of the input units.

Since there is a direct correspondence between a network diagram and its mathematical function, we can develop more general network mappings by considering more complex network diagrams. To be able to run the computation, however, we shall restrict our attention to the case of *feed-forward* networks, in which there are no feed-back loops, as shown in Figure 1 (c).

The activation function of the hidden neurons are usually sigmoidal functions, i.e. $g(a) = \frac{1}{1+e^{-a}}$. The differentiability of this kind of functions leads to a powerful and computationally efficient learning method, called *error backpropagation* [22].

The proposed protocol can deal with sigmoid activation functions and also with linear activation functions. The latter is sometimes desirable because the use of sigmoid units at the outputs would limit the range of possible outputs to the range attainable by the sigmoid.

3 Remote Data Classification

As we already disclosed in the introduction, by remote data classification we mean a scenario where two users, Alice and Bob, remotely cooperate to allow Alice to perform an evaluation or classification applied to her private data set x , by exploiting Bob's knowledge (represented by the classifier own by Bob).

There are two trivial solutions for this problem: Alice can send her data and Bob performs the classification, or Bob sends his classifier and Alice performs the classification. Both these solutions have a clear drawback. Alice doesn't want to disclose to Bob her data nor the output of the classification (in the case of the previous medical example, Alice doesn't want to reveal if she is actually healthy or not). On the other hand, Bob might not be happy to disclose the classifier, since he probably invested time and money in training his network.

Again we can find a trivial solution to these security requirements by resorting to a trusted third party (TTP) that takes inputs from Alice and Bob, and gives back to Alice the output of the computation. Of course this scenario is not the ordinary one: if Alice and Bob distrust each other, it could be hard to agree on a common TTP.

The goal of this work is to propose a protocol for remote data classification that respects the above privacy constraints without resorting to a TTP, and where cryptography and some randomization will play the role of the TTP. As many other protocols in literature, we will make the assumption that both Alice and Bob are semi-honest. This means that they will follow the protocol properly, but they can later analyze the protocol transcripts trying to discover information about other party inputs. We will also assume that Alice and Bob can communicate over a secure (private and authenticated) channel, that can be implemented in practice as a standard SSL connection.

To grant protection against malicious adversary (i.e. adversary that can behave differently from what they are supposed to do), there are standard constructions to produce a protocol secure in the malicious scenario from a semi-honest one [11]. Moreover in our protocol, there is no evidence that a cheating player can discover significant information about other party inputs, and therefore we can assume that semi-honest behavior is forced by external factors i.e. Bob being a service provider that doesn't want to lose his reputation, and Alice is interested in having her data correctly classified.

3.1 Protocol Requirements

In the following, a discussion of the requirements satisfied by the designed protocol to deal with the above scenario is given.

3.1.1 Correctness.

It means that Alice wants to have a correct classification of her data. As Bob gets no output, only Alice has interest in this property. We are not concerned here with the accuracy of Bob's NN, nor we could be. What we mean here is that Bob could maliciously give an incorrect answer to Alice, i.e. in the medical example, a corrupted Bob could be instructed to reply "ill" to a given list of people. We could avoid this kind of cheating by running our protocol over an anonymized network as is [9]. In general as our protocol is only secure against semi-honest adversaries, we can't ensure correctness: to do so, we should transform our protocol into a protocol secure against malicious adversary using standard compiler techniques [11].

3.1.2 Alice's Privacy.

Alice's data are completely protected. In fact, Alice gives her input in an encrypted format and receives the output in an encrypted format. So Alice's privacy relies on the security of the underlying cryptosystem that is, in our case, semantic security.

3.1.3 Bob's Privacy.

While it's clear what we mean with Alice's privacy, it's not the same with Bob's privacy. We need again to refer to the TTP ideal scenario: as long as Alice gets her data classified, she's learning something about Bob classifier. She learns in fact how Bob classifies a vector data. Suppose that she is allowed to run the protocol an arbitrarily number of times: she could submit every (or a very large database of) instances to be classified by Bob, so that she could easily build a table with the classification of every vector, obtaining in that way a classifier that works the same of Bob's one for the classified vector, and that will classify also new vectors.

This result does not depend on the protocol security (we are talking about the TTP ideal scenario) but simply on the fact that from the input and the output of a computation it is always possible to learn something about the computed function (and this is exactly the goal of machine learning). A question is therefore: is it possible to model this kind of attack, where a user interacts many times with the classifier to understand his structure? How many interactions does the attacker need? This question is partially answered in watermarking literature under the name *sensitivity attack*. For a linear classifier (perceptron) the number of iterations to completely disclose it is very low [15]. For a more complex function, like a multi-layer network, it's harder to say, but there are still solutions allowing to extract a local boundary even without knowing anything about the nature of the classifier [7].

The proposed protocol is able to protect the structure and the hidden neurons output of the Bob's NN.

Concerning the structure of the Network, the protocol is designed to completely protect the way the neurons are connected, and to partially protect the number of neurons actually present in the network, so that Alice will get an upper bound B for the number of neurons. Bob can adjust this bound B in a way to achieve a good level of privacy and at the same time an efficient protocol.

Concerning the hidden neurons output, we have to protect them, since they represent an unneeded leakage of information. Given a hidden neuron with sigmoid activation function, it is possible to split its output in two parts, one that we can perfectly protect, while the other will be partially disclosed.

- *State of the Neuron*: this is the most important part of the output. Depending on the sign of the weighted sum of the inputs, every neuron can be *on* (i.e. output > 0.5) or *off* (i.e. output < 0.5). We will perfectly protect this information, flipping it with probability one half, achieving a one-time pad kind of security.
- *Associated Magnitude*: the activation of every neuron has also a magnitude, that gives information about "how much" the neuron is *on* or *off*. We will hide those values in a large set of random values, in such a way that Alice will not learn which values are actually part of the computation and which ones are just random junk. Of course the bigger the set is, the more privacy it comes, and more computational resources are needed.
- *Neuron's Position*: the last trick to achieve security is to completely randomize the position of the hidden neurons in such a way that Alice will not discover which outputs correspond to which neurons. Therefore Alice may learn a bit of information at every execution (meaning something about the magnitudes of them), but she'll not be able to correlate those information in order to gain advantage from repeated executions of the protocol.

3.1.4 Round Complexity.

We would like to run the computation with no interaction (except the minimum needed to input the vector and get the output). Unluckily, there are only few kinds of algorithms that can be computed in such a way, that is NC^1 circuits [23], and $NLOGSPACE$ problem [4]. Round complexity, i.e. the number of messages between the players, is an important measure for the complexity of a cryptographic protocol. The round complexity of our protocol is given by the number

of layers of the network. Given that a network with three layers can approximate any function, we can consider our protocol to have a constant round complexity.

4 Privacy-Preserving Protocol for Remote Data Classification

We will continue to use the notation introduced before, together with some new symbols to refer to the encrypted version of the data. The input of the neuron is \mathbf{x} and its encrypted version is \mathbf{c} , while the encrypted version of the activation a of every neuron will be referred as d .

4.1 Building blocks

We describe first of all the building blocks that are required for the correct behavior of the proposed protocol.

4.1.1 Homomorphic Encryption.

The chosen public-key cryptosystem to instantiate our protocol is the Damgård-Jurik modification [8] of the Paillier encryption scheme [20].

This cryptosystem is based on the hardness to decide the n -th residuosity of elements modulo n^{s+1} , where n is an RSA modulo. At the end, the encryption and the decryption procedures are the following:

Set-up: select p, q big primes. Let $n = pq$ be the public key, while the secret key, called λ , is the least common divisor between $(p-1)$ and $(q-1)$.

Encryption: let $m \in \mathbb{Z}$ be the plaintext, and s such that $n^s > m$. Select a random value $r \in \mathbb{Z}_n^*$; the encryption c of m is:

$$c = E_{pk}(m, r) = (1 + n)^m r^{m^s} \pmod{n^{s+1}}$$

Decryption: the decryption function D_{sk} depends only on the ciphertext, and there is no need to know the random r in the decryption phase. We refer to the original paper for the complete description.

The main advantage of this cryptosystem is that the only parameter to be fixed is n , while s can be adjusted according to the plaintext. In other words, unlike other cryptosystems, where one has to choose the plaintext m to be less than n , here one can choose an m of arbitrary size, and then adjust s to have $n^s > m$ and the only requirement for n is that it must be unfeasible to find its factorization.

The trade-off between security and arithmetic precision is a crucial issue in secure signal processing applications. As we will describe in more detail in the Appendix, the possibility to quantize the input samples and to work with an arbitrary precision allows us to neglect that we are dealing with integer modular numbers, so that we can consider to have an arbitrarily accurate non-integer homomorphic encryption scheme.

For the sake of simplicity from now on we will indicate the encryption just as $c = E(m)$, as the keys are chosen once and are the same for all the protocol length, and the random parameter r is just to be chosen at random. If $x_1 = x_2$ we will write $E(x_1) \sim E(x_2)$. The encryption of a vector $\mathbf{c} = E(\mathbf{x})$ will be simply the vector composed of the encryption of every component of the plaintext vector.

As said, this cryptosystem satisfies the homomorphic property so, given two plaintexts m_1 and m_2 , the following equalities are satisfied:

$$D(E(m_1) \cdot E(m_2)) = m_1 + m_2 \tag{1}$$

and

$$D(E(m)^a) = am. \tag{2}$$

Concerning the security of this cryptosystem, it provides semantic security against chosen-plaintext attacks; more details can be found in the original paper [20].

4.1.2 Privacy Preserving Scalar Product (PPSP).

A secure protocol for the scalar product allows Bob to compute an encrypted version of the scalar product between an encrypted vector given by Alice $\mathbf{c} = E(\mathbf{x})$, and one vector \mathbf{w} owned by Bob. The protocol guarantees that Bob gets nothing, while Alice gets an encrypted version of the scalar product that she can decrypt with her private key. Such a protocol is easily achievable exploiting the two homomorphic properties shown in the Equations 1 and 2:

Input: $\mathbf{c} = E(\mathbf{x})$; Bob: \mathbf{w}
Output: Alice: $d = E(\mathbf{x}^T \mathbf{w})$
 PSPP($\mathbf{c}; \mathbf{w}$)
 (1) Bob computes $d = \prod_{i=1}^N c_i^{w_i}$
 (2) Bob sends d to Alice

After receiving d , Alice can decrypt this value with her private key to obtain the weighted sum a .

It is worth observing that though the above protocol is a secure one in a cryptographic sense, some knowledge about Bob's secrets is implicitly leaked through the output of the protocol itself. If, for instance, Alice can interact N times with Bob (where $N = |\mathbf{x}| = |\mathbf{w}|$ is the size of the input vectors), she can completely find out Bob's vector, by simply setting the input of the i -th iteration as the vector with all 0's and a 1 in the i -th position, for $i = 1, \dots, N$. If we use the scalar product protocol described above to build more sophisticated protocols, we must be aware of this leakage of information. This is again a *sensitivity attack*, as introduced before. Note that the problems stemming from sensitivity attacks are often neglected in the privacy preserving computing literature.

4.1.3 Evaluation of the Activation Function

If the function g is known and it's invertible, like in the case of the sigmoid function, the information given by a or $y = g(a)$ is the same. So Bob can simply give a to Alice that can compute y by herself.

4.2 Perceptron Protocol

The blocks described above allow to run a privacy preserving remote data classification protocol for a single layer network. Let us suppose that the network has I inputs and 1 output. The protocol is the following:

Input: $\mathbf{c} = E(\mathbf{x})$; Bob: \mathbf{w}
Output: Alice: classification of \mathbf{x}
 PERCEPTRON($\mathbf{c}; \mathbf{w}$)
 (1) Alice and Bob run the PPSP protocol.
 (2) Alice decrypts the output $a = D(d)$
 (3) Alice computes $g(a)$

If the network has more than one output neuron, say O , just run in parallel O instances of the previous protocol.

4.3 Handling with Hidden Neurons

We consider now the more interesting case of a feedforward network. As already defined, a feedforward network is composed by N neurons that can be ordered in a way that neuron j gets in input the output of a finite set I_j of neurons

having index lower than j , like the ones in Figure 1. We use this ordering to label every neuron. The weight of the connection from neuron i to neuron j is indicated by w_{ij} . The input vector of neuron j is \mathbf{x}_j while the associated weights vector \mathbf{w}_j . The aim is now to protect the output of hidden neurons and the network topology.

4.3.1 Hidden Neurons Output Protection.

In the perceptron protocol, Bob gives to Alice the linear sum a of the output neurons, and then Alice computes by herself the activation function output. The simple iteration of the perceptron protocol for every neuron in the network will disclose the activation value of every hidden neuron, but Alice is not supposed to get this information.

The linear sum a can be viewed as $a = \text{sign}(a) \cdot |a|$. Depending on $\text{sign}(a)$ the output of the sigmoid function will be “almost 1” (i.e. $0.5 \leq y < 1$) or “almost 0” (i.e. $0 < y \leq 0.5$). We can perfectly protect the value $\text{sign}(a)$ by exploiting the fact that the sigmoid function is actually antisymmetric as shown in Figure 2, i.e. $g(-a) = 1 - g(a)$.

Bob can randomly change the sign of a just before sending it to Alice thanks to the homomorphic property, since $E(a)^{-1} = E(-a)$. Then Alice can decrypt as usual the value received, compute the activation function on the received input $g(-a)$ and send $E(g(-a))$ back to Bob. Bob can recover the value he needs simply performing the subtraction (only for the scrambled values), that is $E(g(a)) = E(1 - g(-a)) = E(1)E(g(a))^{-1}$.

In this way Alice will not discover which (nor how many) neurons are actually activated or not. We will deal with the protection of the value $|a|$ later.

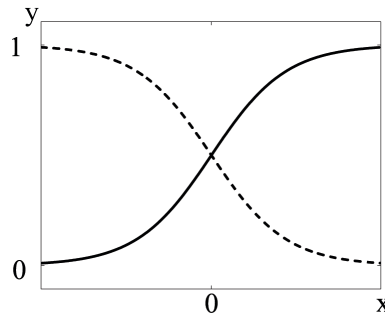


Fig. 2: The Sigmoid is antisymmetric with respect to $(0, 1/2)$, i.e. $g(-a) = 1 - g(a)$.

4.3.2 Network Embedding.

To protect the network topology, i.e. the number of hidden neurons N_H and the way they are connected, we will embed the network in a multilayer network, composed of L layers of M neurons each. Of course $LM \geq N_H$. The added $LM - N_H$ neurons will be called *fake neurons*. They have to look the same as the real neurons and so they will be initialized with some incoming connections from other neurons, with random weights. They will not influence the computation as no real neurons will take their outputs as input.

A good embedding is one where every neuron only gets inputs from neurons that are in previous layers. An example of a legal embedding of the network in Figure 1 (b) is given in Figure 3.

In this way Alice will only learn an upper bound LM for the actual number of hidden neurons N_H , while she will not learn the way they are connected or the connection weights, as Bob can perform the PPSP for every neuron by himself just picking up the necessary inputs and make the weighted sum with his private weights. The number L also gives a bound about the longest path between input and output neurons. Instead M is not the upper bound of the incoming connection for a neuron, given that we can split one original layer into two or more layers in the embedded network. Every neuron in layer l can take inputs from any neuron belonging to any previous layer.

The more we increase L , the more we protect the network topology. At the same time L will be the number of interaction round of the protocol, so we have to find a tradeoff between round complexity and network topology protection.

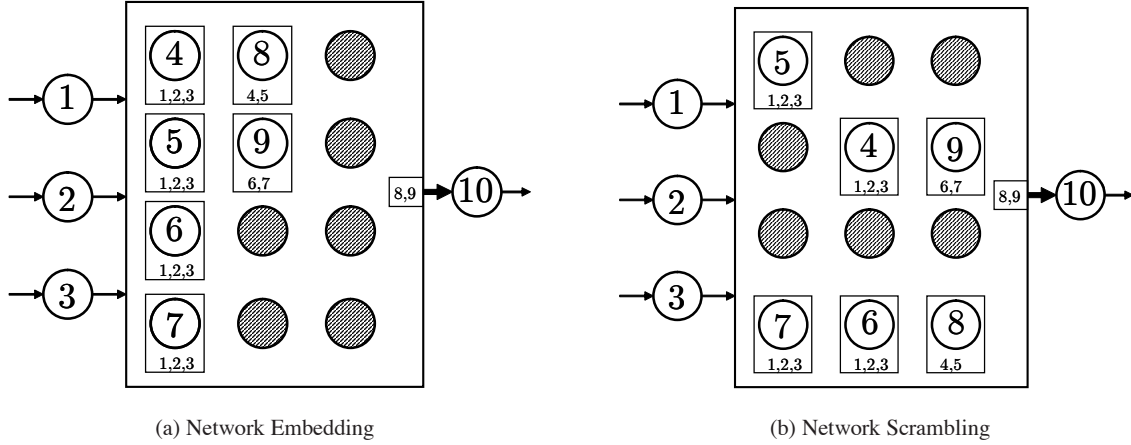


Fig. 3: Figure (a) is a natural embedding of the network in Figure 1 (b) into a 3×4 network. The big box refers to the protected zone. The numbers under the neurons indicate which inputs are given to that neuron. Filled neurons refer to fake neurons, that are meaningless for the actual computation. In the setup phase we need to assign some inbound connection for these neurons, to be sure that their output will be undistinguishable from that of real hidden neurons. Figure (b) shows a legal scrambling of the network.

4.3.3 Network Randomization.

At this point we have to deal with the protection of the value $|a|$. We are concerned with the disclosing of $|a|$ because if Alice is allowed to run the protocol several times, she can use this additional information to actually understand the network weights. The solution we propose is to scramble all the network at every different execution. A legal scrambling is one that preserves the property of the embedding, i.e. every neuron only gets input from neurons belonging to previous layers, as the one shown in Figure 3 (b). We note that there are at least $L \cdot M!$ legal scramblings, the ones that permute only neurons between the same layer. Increasing the number of neurons per layer M we can get a higher number of different permutations and so a better protection for the network, at the cost of a higher computational burden.

We define the *embedding ratio* as the ratio between the number of hidden neurons of the embedded network and the number of hidden neurons in the original network LM/N_H . As this ratio increases, Alice will see more meaningless values, and therefore she won't be able to understand which values refers to real hidden neurons and which ones don't. Together with the network scrambling, this is a countermeasure to Alice's attempt to run a sensitivity attack against the protocol.

4.4 Multi-Layer Network Protocol

The final protocol is the following:

System Setup: Bob chooses L, M and finds a legal embedding of his network in the $L \times M$ one.

Execution Setup: Alice and Bob agree on a parameter s required by the used cryptosystem, and on a quantization factor Q , that will be used to obtain integer values from the input samples. Alice generates a pair of public and

private keys and gives the public one to Bob. Bob will randomly scramble the neuron positions in the network to reach another legal configuration.

Execution: we will consider real neurons and fake ones as the same, as there's no difference, but the fake one's output will never be used again.

Input: Alice: \mathbf{c} ; Bob: $\mathbf{w}_j^{(i)}$ with $i = 1, \dots, L, j = 1, \dots, M$
Output: Alice: classification of x
 PPDC($\mathbf{c}; \{\mathbf{w}_j^{(i)}\}_{i=1, \dots, L; j=1, \dots, M}$)

- (1) **for** $i = 1$ **to** L
- (2) **for** $j = 1$ **to** M
- (3) Bob runs PPSP($\mathbf{c}_j^{(i)}; \mathbf{w}_j^{(i)}$) and gets $d = E(a)$
- (4) Bob picks $t_j \in \{+1, -1\}$ at random
- (5) **if** $t_j = -1$
- (6) Bob sets $d = d^{-1}$
- (7) Bob sends d to Alice
- (8) Alice decrypts d , evaluates $g(a)$, and sends $z = E(g(a))$ back to Bob
- (9) **if** $t_j = -1$
- (10) Bob sets $z = E(1)z^{-1}$
- (11) **for** $j = 1$ **to** O //out-degree of the network
- (12) Bob runs PPSP($\mathbf{c}_j; \mathbf{w}_j$) and gets $d = E(a)$
- (13) Bob sends d to Alice
- (14) Alice decrypts d and evaluates her output $g(a)$

The security of this protocol follows from the previous considerations.

5 Implementation of the Protocol

In this section a practical implementation of the proposed protocol is described, and a case study execution that will give us some numerical results in terms of computational and bandwidth resources needed is analyzed.

5.0.1 Client-Server Application.

We developed a Java application based on Remote Method Invocation technology³. The software, which makes use of a modified implementation of the Damgård-Jurik cryptosystem available on Jurik's homepage⁴, is composed of two parts: the client and the server. The former sets the environment creating a couple of public/private keys and choosing the number of bits of the modulus, and it chooses a server to connect to. The latter can load several neural networks into the system and choose an appropriate quantization factor and embedding ratio.

5.0.2 Experimental Data.

Two datasets were selected from the UCI Machine Learning Repository [21], and two kinds of neural networks were trained starting from those data sets:

- *Sonar*: this dataset refers to the classification of sonar signals by means of a neural network; the dataset contains patterns corresponding to sonar signals bounced off a metal cylinder (bombs) and signals bounced off a roughly cylindrical rock; we have trained a NN with the standard backpropagation algorithm, containing 60 input neurons, 12 hidden neurons, and 1 output neuron.

³ <http://java.sun.com/javase/technologies/core/basic/rmi/>

⁴ <http://www.daimi.au.dk/~jurik/research.html>

- *Nursery*: this dataset was derived from a hierarchical decision model originally developed to rank applications for nursery schools [12]; we have trained a NN with 8 input neurons, 20 hidden neurons, and 5 output neurons.

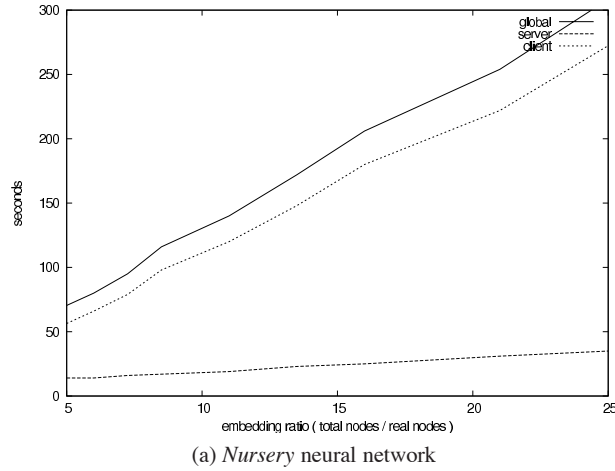
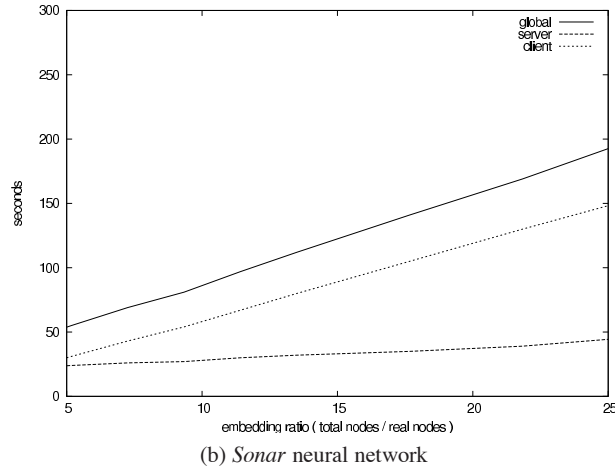
(a) *Nursery* neural network(b) *Sonar* neural network

Fig. 4: Protocol execution time, according to different values of the embedding ratio.

5.0.3 Experimental Setup.

We loaded these NNs into the server changing the embedding ratio at every execution. The quantization factor has been set to $Q = 10^9$, obtaining a reasonable accuracy in the computation to prevent quantization error. In the key generation algorithm the key length has been set to $n = 1024$ to obtain a reasonable level of security [2].

Then we deployed the application on two mid-level notebooks, connected on a LAN. The execution time increases linearly with respect to the embedding ratio, as also the communication overhead does. Choosing carefully the embedding ratio, we can find the desired trade off between execution time and achieved security. Experiments showed that an embedding ratio value of 10 offers a reasonable execution time. Of course the level of security we need is mostly application-based. Results are pictorially represented in Figure 4, where we can distinguish between the running time on the client, on the server, and the total time. Given that the application was executed on a LAN, the communication time is negligible. The amount of exchanged bytes is reported in Table 1 for some values of the embedding ratio. It

is worthy to note that even if the *Sonar* NN has 73 neurons (60 input + 12 hidden + 1 output) while the *Nursery* NN has only 33 neurons (8 input + 20 hidden + 5 output), the former is computed in shorter time. This is due to the fact that the only piece of the network that has to be embedded are the hidden neurons⁵. Therefore the *Nursery* NN, having more hidden neurons than the *Sonar* NN, needs more time and bandwidth for the computation.

Neural Network	embedding ratio				
	5	10	15	20	25
sonar	153kB	256kB	359kB	470kB	579kB
nursery	232kB	368kB	544kB	722kB	894kB

Table 1: bandwidth occupation

5.0.4 Workload.

Looking at client and server side workload we can notice that something (apparently) strange happens: the client side workload is greater than the server side one, despite the server seems to do more work than anyone else.

Actually, taking into consideration the protocol, during the evaluation of every hidden layer the client performs encryptions and decryptions once for every neuron while the server has to do a lot of multiplications and exponentiations at every step (that is, for every neuron) to compute the activation value. This is why we are allowed to think that the server side workload has to be greater than the client side one.

Nevertheless, this is not true. We should consider also the nature of things, not only the amount of them. In fact, the client performs heavier operations than the server (encryptions and decryptions are more expensive than multiplications and exponentiations, of course) and so the client side workload is greater than the other one. So, as shown in Fig. 4, with a little amount of nodes the global workload seems to be fairly shared among the parts, and increasing the number of hidden nodes the client side workload grows faster and the performance gets worse.

That behavior allows the protocol to work fine with NNs that are composed by a relatively small set of hidden neurons (that is, between 100 and 300 hidden neurons) because of the nature of the operation involved, no matter which kind of computer the server offers. In any case, we can notice in Fig. 4 that also with an embedding ratio equal to 25 it is possible to achieve a satisfactory computation time, so that it is possible to handle without problems a NN with more than 500 hidden neurons, as we did with sonar and nursery neural networks during the experimental tests.

As far as we know this is the first protocol of this kind that was implemented, as no experimental results are mentioned in [6]. Therefore we can't compare our approaches.

6 Conclusions

Several modern artificial intelligence applications require the protection of the privacy of the data owner, unwilling to reveal his/her input data to the owner of the classifier. In this framework, the availability of tools to process data and signals directly in the encrypted domain allows to build secure and privacy preserving protocols solving the mentioned problem. Given the central role that neural network computing plays in artificial intelligence field, a protocol for NN-based privacy-preserving computation has been designed, where the knowledge embedded in the NN as well as the data the NN operates on are protected. The proposed protocol relies on homomorphic encryption for the linear computations, and on a limited amount of interaction between the NN owner and the user for non linear operations. However, the interaction is kept to a minimum, without resorting to general multiparty computation protocols. Any

⁵ The input and output layers represent the input and the output of the function: adding fake inputs or outputs, or scrambling their position will result in a meaningless computation.

unnecessary disclosure of information has been avoided, protecting all the internal computations, so that at the end the user will only learn the final output of the NN computation.

Acknowledgements The work described in this paper has been supported in part by the European Commission through the IST Programme under Contract no 034238 - SPEED. The information in this document reflects only the authors views, is provided as is and no guarantee or warranty is given that the information is fit for any particular purpose. The user thereof uses the information at its sole risk and liability.

References

1. Agrawal, R., Srikant, R.: Privacy-preserving data mining. In: SIGMOD '00: Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data, pp. 439–450 (2000)
2. Barker, E., Barker, W., Burr, W., Polk, W., Smid, M.: Recommendation for Key Management—Part 1: General. NIST Special Publication pp. 800–57 (2005)
3. Barni, M., Orlandi, C., Piva, A.: A privacy-preserving protocol for neural-network-based computation. In: MM&Sec '06: Proceeding of the 8th Workshop on Multimedia and Security, pp. 146–151 (2006)
4. Beaver, D.: Minimal-latency secure function evaluation. Proc. of EUROCRYPT00, LNCS pp. 335–350 (2000)
5. Bishop, C.: Neural Networks for Pattern Recognition. Oxford University Press (1995)
6. Chang, Y., Lu, C.: Oblivious polynomial evaluation and oblivious neural learning. Theoretical Computer Science **341**(1), 39–54 (2005)
7. Comesana, P., Perez-Freire, L., Perez-Gonzalez, F.: Blind newton sensitivity attack. Information Security, IEE Proceedings **153**(3), 115–125 (2006)
8. Damgård, I., Jurik, M.: A generalisation, a simplification and some applications of Paillier's probabilistic public-key system. In: Public Key Cryptography, pp. 119–136 (2001)
9. Dingedine, R., Mathewson, N., Syverson, P.: Tor: The second-generation onion router. In: Proceedings of the 13th USENIX Security Symposium (2004)
10. Fouque, P., Stern, J., Wackers, G.: CryptoComputing with rationals. Financial-Cryptography.-6th-International-Conference, Lecture-Notes-in-Computer-Science **2357**, 136–46 (2003)
11. Goldreich, O., Micali, S., Wigderson, A.: How to play any mental game or a completeness theorem for protocols with honest majority. In: STOC '87: Proceedings of the twentieth annual ACM symposium on Theory of computing, pp. 218–229 (1987)
12. Gorman, P., et al.: UCI machine learning repository (1988). URL <http://archive.ics.uci.edu/ml/>
13. Hornik, K., Stinchcombe, M., White, H.: Multilayer feedforward networks are universal approximators. Neural Netw. **2**(5), 359–366 (1989)
14. Jha, S., Kruger, L., McDaniel, P.: Privacy Preserving Clustering. 10th ESORICS (2005)
15. Kalker, T., Linnartz, J.P.M.G., van Dijk, M.: Watermark estimation through detector analysis. In: ICIP98, vol. I, pp. 425–429. Chicago, IL, USA (1998)
16. Kantarcioglu, M., Vaidya, J.: Privacy preserving naive bayes classifier for horizontally partitioned data. In: In IEEE Workshop on Privacy Preserving Data Mining (2003)
17. Lapedes, A., Farber, R.: How neural nets work. Neural information processing systems pp. 442–456 (1988)
18. Laur, S., Lipmaa, H., Mielikäinen, T.: Cryptographically private support vector machines. In: KDD '06: Proceedings of the 12th ACM SIGKDD International Conference on Knowledge discovery and data mining, pp. 618–624 (2006)
19. Lindell, Y., Pinkas, B.: Privacy preserving data mining. In Advances in Cryptology - CRYPTO '00, volume 1880 of Lecture Notes in Computer Science **1880**, 36–54 (2000)
20. Pailler, P.: Public-key cryptosystems based on composite degree residuosity classes. In: Proceedings of Eurocrypt'99, Lecture Notes in Computer Science vol. 1592, pp. 223–238. Springer-Verlag (1999)
21. Rajkovic, V., et al.: UCI machine learning repository (1997). URL <http://archive.ics.uci.edu/ml/>
22. Rumelhart, D., Hinton, G., Williams, R.: Learning internal representations by error propagation. Parallel distributed processing: Explorations in the microstructure of cognition **1**, 318–362 (1986)
23. Sander, T., Young, A., Yung, M.: Non-Interactive CryptoComputing For NC 1. IEEE Symposium on Foundations of Computer Science pp. 554–567 (1999)
24. Yang, Z., Wright, R.N.: Improved privacy-preserving bayesian network parameter learning on vertically partitioned data. In: ICDEW '05: Proceedings of the 21st International Conference on Data Engineering Workshops, p. 1196 (2005)
25. Yao, A.: How to generate and exchange secrets. In: IEEE FOCS'86 - Foundations of Computer Science, pp. 162–167 (1986)

Appendix: Handling Non-integer Value

In this appendix we will describe how we can obtain an arbitrarily accurate real number Damgård-Jurik cryptosystem. First of all we map, as usual, the positive numbers in $\{0, \dots, \frac{n^s-1}{2}\}$, and the negative ones in $\{\frac{n^s-1}{2} + 1, \dots, n^s - 1\}$,

with $-1 = n^s - 1$. Then, given a real value $x \in \mathbb{R}$, we quantize it with a quantization factor Q , and approximate it as $\bar{x} = \left\lfloor \frac{x}{Q} \right\rfloor \simeq \frac{x}{Q}$ for a sufficiently thin quantization factor. Clearly the first homomorphic property still stands i.e.

$$D(E(\bar{x}_1) \cdot E(\bar{x}_2)) = \bar{x}_1 + \bar{x}_2 \simeq \frac{x_1 + x_2}{Q}.$$

This allows Bob to perform an arbitrarily number of sums among ciphertexts. Also the second property holds, but with a drawback. In fact:

$$D(E(\bar{x})^a) = a \cdot \bar{x} \simeq \frac{a \cdot x}{Q^2}$$

The presence of the Q^2 factor has two important consequences: (i) the size of the encrypted numbers grows exponentially with the number of multiplications; (ii) Bob must disclose to Alice the number of multiplications, so that Alice can compensate for the presence of the Q^2 factor. The first drawback is addressed with the property of Damgård-Jurik cryptosystem that allows us, by increasing s , to cipher bigger numbers. The second one imposes a limit on the kind of secure computation that we can perform using the techniques proposed here. Luckily in our application we perform only one multiplication for each ciphertext in the scalar product protocol.

An estimation of the s parameter to be chosen in a way that the value inside the ciphertext after the computation will fit into n^s is the following: let X be the upper bound for the norm of Alice's input vector, and W an upper bound for the weight vectors norm. Every scalar product computed in the protocol is then bounded by $|\mathbf{x}| |\mathbf{w}| \cos(\mathbf{x}\mathbf{w}) \leq XW$. Given a modulo n sufficiently large for security purposes, it is possible to select the parameter s such that:

$$s \geq \left\lceil \log_n \frac{2XW}{Q^2} \right\rceil$$

(the factor 2 is due to the presence of both positive and negative values). Other solutions for working with non integer values can be found in [6] where a protocol to evaluate a polynomial on floating-point numbers is defined (but the exponent must be chosen in advance), and [10], where a cryptosystem based on lattice properties allowing computation with rational values is presented (even in this case, however, a bound exists on the number of multiplications that can be carried out to allow a correct decryption).