



A Quantile-Based Watermarking Approach for Distortion Minimization

Maikel Lázaro Pérez Gort^(✉), Martina Olliaro, and Agostino Cortesi

Ca' Foscari University of Venice, DAIS, Via Torino 155, 30172 Mestre, Venice, Italy
{maikel.perezgort,martina.olliaro,cortesi}@unive.it

Abstract. Distortion-based watermarking techniques embed the watermark by performing tolerable changes in the digital assets being protected. For relational data, mark insertion can be performed over the different data types of the database relations' attributes. An important goal for distortion-based approaches is to minimize as much as possible the changes that the watermark embedding provokes into data, preserving their usability, watermark robustness, and capacity. This paper proposes a quantile-based watermarking technique for numerical cover type focused on preserving the distribution of attributes used as mark carriers. The experiments performed to validate our proposal show a significant distortion reduction compared to traditional approaches while maintaining watermark capacity levels. Also, positive achievements regarding robustness are visible, evidencing our technique's resilience against subset attacks.

Keywords: Distortion reduction · Numeric distribution · Quantile · Robust watermarking · Watermark capacity

1 Introduction

With the easy access and spreading of digital content through the Internet, data copyright protection faces more and more challenges every day. Digital watermarking has become a handy tool to deal with false ownership claims and illegal data copy distribution. The general idea of watermarking techniques consists of adding hidden content (i.e., the watermark) into the protected data. Under demands, watermarks can be extracted and used as evidence of rightful ownership and data tampering, among others. Considering that watermarking is not based on blocking access or copying data, their portability benefits (e.g., allowing data to reach the target communities) are never affected. For the sake of authenticity and trust, usability and intellectual property of data must be protected at all costs.

According to the distortion criterion, watermarking techniques can be classified as distortion-free or distortion-based [2, 16]. Distortion-free techniques generate the watermark from a particular digital asset copy (or embed it into the

data without performing updates) [14, 17]. In contrast, distortion-based techniques perform watermark embedding by modifying the data as long as changes are permissible and do not compromise their usability [9].

Distortion-based watermarking techniques are characterized by two main processes: (i) watermark embedding, (ii) and its extraction. The embedding process first encodes the watermark and then performs its injection into the data. If the encoding uses a meaningful source (e.g., an image file, an audio stream, or a text document) for watermark generation, the watermark is classified as meaningful. Otherwise, it is classified as meaningless [7]. Instead, the extraction process detects every mark from the data and then carries out their extraction to proceed with the watermark reconstruction. Some techniques only perform the detection phase, stating the presence or absence of the watermark in the data [4]. Performing both embedding and extraction processes requires at least one parameter defined as the Secret Key. This parameter must remain secret, and it has to keep the same value for both processes [1].

In most cases, distortion-based approaches are oriented to ownership protection and must be resilient against attacks focused on compromising watermark detection. For this reason, they are classified as robust techniques.

One of the major challenges for distortion-based techniques is guaranteeing data usability despite the changes performed on them. This is hard to achieve considering that according to the robustness requirement, a significant number of marks must be inserted into the data to allow the watermark signal persistence despite attacks. Then, the higher the number of marks inserted, the higher the distortion over the data. Thus, the number of marks embedded into the digital assets (defined as watermark capacity) is inversely proportional to the watermark imperceptibility in the data. Indeed, the imperceptibility requirement is expected to be accomplished as long as the distortion does not cause degradation of data usability.

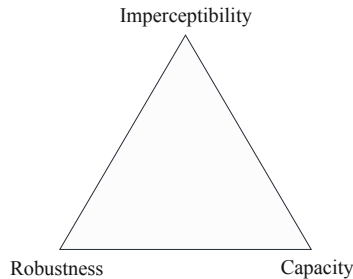


Fig. 1. Trade-off among robustness, imperceptibility, and capacity requirements [10].

There is a trade-off that watermarking techniques must deal with regarding robustness, capacity, and imperceptibility requirements (see Fig. 1). The strong link among them and the equality of their relevance for the technique's success is represented as an equilateral triangle. As long as one of them is affected, the

others will be impacted for better or worst. For example, a common approach for distortion reduction is to reduce the watermark capacity, negatively affecting the technique’s robustness. Indeed, it is not possible to significantly increase the imperceptibility without having a negative influence over robustness.

1.1 Paper Contribution

In this paper, we propose a strategy to benefit watermark imperceptibility in techniques embedding marks in numerical attributes (a.k.a., numerical cover type watermarks) of database relations, without affecting watermark capacity.

Our main goal is to preserve the numerical distribution of the columns used as carriers as much as possible, avoiding some values from moving from quantiles defined to control the distribution. When the new value containing the mark changes quantile after the embedding, the marking should not be rolled back since this would reduce the watermark capacity. Instead of allowing values changes between quantiles, we propose a mechanism for performing mark embedding in other carriers’ distributions allowed regions.

The experiments performed show a significant enhancement of imperceptibility once numerical distribution is kept as similar as possible with respect to the original unwatermarked columns. We used scatter statistical metrics to compare the effects of watermark embedding of our approach vs. conventional embedding. Also, we applied the Kullback-Leibler divergence to measure the relative entropy between the distribution of the original data and the one resulting from the watermark embedding. Since the watermark capacity is not affected, robustness improves, making it more difficult for attackers to compromise watermark signal detection.

1.2 Paper Structure

The rest of the paper is organized as follows. Section 2 offers details of the theoretical background, presenting commonly used notations in the relational data watermarking research field. Also, in this section, the related work (mostly focused on approaches oriented to distortion-reduction) is given. Section 3 presents our proposal, depicting the benefits and downsides of each strategy of quantile-based numerical distribution preservation. Section 4 presents the experimental results, mainly oriented to show the behavior of robustness, capacity, and imperceptibility watermark requirements. Section 5 concludes.

2 Theoretical Background

Contrary to multimedia data, effects of the watermark (WM) embedding into relational data are not perceived directly by human systems (e.g., human visual system, human auditory system). Instead, a Middle Coded-based Layer (MCL) composing management information systems processes the data and delivers it to users in more suitable formats such as digital reports. This has an important

consequence. Indeed, WM imperceptibility does not depend on human systems limitations but on the processes implemented by MCL, which are often based on business rules. Following that principle, it may appear that WM capacity benefits from the inability of direct human perception over relational data changes. Nevertheless, as long as digital systems generate outputs from the watermarked data (having others using them as inputs), the slightest changes will drive drastic consequences.

Among their classification criteria, relational data watermarking defines the technique type according to the data type of attribute selected in the relation R to perform the WM embedding (also known as mark carriers). Some techniques use textual attributes, being classified as textual cover type approaches (e.g., Al-Haj & Odeh [3], Pérez Gort et al. [6]). Others are focused on numerical cover types (e.g., Rani et al. [12], Hou & Xian [8], Zhao et al. [18]), etc. For numerical cover type approaches, it is very common to perform WM embedding by inserting each mark in one position selected from a given range of less significant bits (*lsb*) of the carrier attribute numerical value binary representation.

Even if just the first *lsb* is changed, the impact at column level could be higher compared to at attribute-value level. Also, depending on the MCL implemented processes, changes might not be tolerable if a general description of the behavior of the data is used for decision making. Some changes at single-value level might appear tolerated, but the effects over the whole set of data might contradict database purposes.

2.1 Related Work

In 2002, Agrawal & Kiernan [2] highlighted for the first time the need for watermarking relational data for ownership protection and formalized the so-called AHK watermarking algorithm. Precisely, based on the condition that some attribute's values can tolerate changes (as long as data usability is preserved), they proposed to mark only numeric columns. Embedding is performed at bit level, where carriers are pseudo-randomly selected according to a Secret Key (SK). However, this technique has proved to be vulnerable to simple attacks (e.g., bit flipping and updates attacks) due to the meaningless of WM information (i.e., bit pattern). Usability control is based on the number of *lsb* available for marking in an attribute and the number of marked tuples, while constraints deployed over the database are ignored.

Statistic metrics describing the numerical distribution featuring the attribute selected for WM embedding are a good reference to appreciate the general changes performed compared to the distribution before the embedding.

In 2004, Sion et al. [15] proposed a numerical cover type technique performing embedding of marks at bit level. For this case, usability maintenance is done by data statistics preservation. Also, the marking of selected tuples is performed according to database constraints and an error range allowed for data, using the Mean Squared Error (MSE) as reference. Nevertheless, this proposal requires tuple ordering to define subsets identifying some tuples as group bounds, being vulnerable to subset reverse order, tuple updates, and deletion attacks.

In 2010, Sardroudi & Ibrahim [13] proposed a new watermarking technique using as WM source a binary image. Given a relation, their schema embeds marks only in one numerical attribute, focusing on guaranteeing robustness and minimizing data variation by flipping the first *lsb* depending on the value of the mark embedded. This technique shows good results against subset reverse order attacks. Nevertheless, capacity is often affected by the partial embedding of the watermark, making it vulnerable to other malicious operations such as subset update attacks.

Pérez Gort et al. [11], in 2017, proposed a technique extending Sardroudi & Ibrahim's scheme, where the embedding is performed over more than one attribute per tuple according to one parameter defined as Attribute Fraction (AF). In this case, distortion reduction at bit level is also performed, but flipping all *lsbs* to the right of the one selected for mark embedding, depending on their values and the value of the mark. Nevertheless, reducing distortion at the bit level does not always benefit the numerical distribution of the carrier column. In that sense, WM embedding is performed blindly and the general quality of data could be compromised.

Techniques based on the AHK [2] algorithm select $\omega \approx \eta/\gamma$ tuples to mark out of the η stored in the relation R, being $\gamma \in [1, \eta]$ the Tuples Fraction (TF) representing the inverse of the marking density. For each tuple selected, an attribute (out of v attributes) is chosen and the binary representation of the contained value is used for inserting the mark. Sardroudi & Ibrahim's [13] technique increases the link between the watermark source and R. To this aim, each pixel *pseudo-randomly* selected from the binary image used as WM source is *xored* with one of the most significant bits (*msb*) of a range given as parameter (denoted as β) of the value where the mark will be embedded. Finally, the *lsb* position is selected from a given number of bits available for marking (denoted as ξ), and the mark generated is embedded into it. Considering the approaches just mentioned embed only one mark per tuple, Pérez Gort et al. [11] extends the embedding to more than one attribute by defining AF (denoted as $\delta \in [1, v]$), where $\delta = 1$ forces all attributes of the selected tuples to be marked.

3 Proposed Approach

Note that none of the approaches discussed in the previous section analyzes the distortion caused by WM embedding from a numerical distribution point of view. This is a critical issue since, depending on the distribution variation, data can result useless after the embedding, according to the data owner's goals. In this work, besides taking care of the distortion from the binary level perspective, also different proposals are presented to preserve each attribute's distribution. Our main goal is to maintain as similar as possible the resulting distributions after WM embedding with respect to the one each attribute had before R being distorted.

Formally, let us denote by \mathcal{D}_i the distribution of the attribute i before the WM embedding, and \mathcal{D}'_i after the embedding. If we denote by \equiv the equivalence relation between distributions, we aim to achieve the following condition:

$$\forall i \in [0, v - 1] : \mathcal{D}_i \equiv \mathcal{D}_i' \quad (1)$$

We start by fragmenting each distribution \mathcal{D}_i in g quantiles (see Fig. 2a)) to prevent the distribution from suffering high variations during the WM embedding.

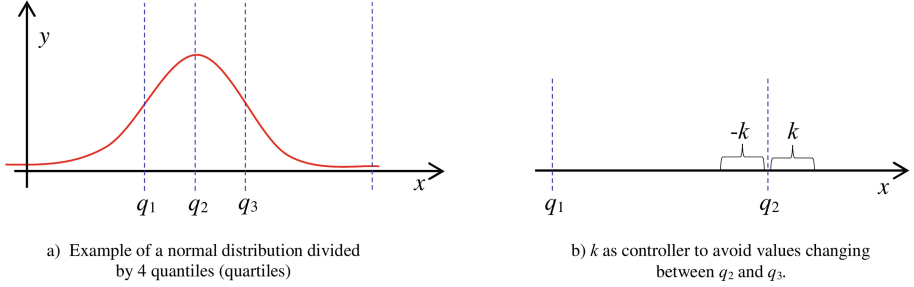


Fig. 2. Numerical distribution split into 4 quantiles (known as quantiles).

Besides the number of quantiles g , also a threshold to control the set of values restricted to be marked in the limits of the quantile (denoted as k) is considered in order to prevent distribution variations (see Fig. 2b)).

The main lines of action followed in this work are: (i) reversing the embedding and preventing the value from being marked, (ii) performing the embedding by assigning values to other distributions as long as quantile changes are not carried out. Each one of these alternatives is detailed below.

3.1 First Action: Mark Embedding Cancellation

Once a value v is selected to be marked, its quantile is located according to $[q_l, q_u] = \mathcal{Q}(\mathcal{D}_i, v, g)$, where \mathcal{Q} is the function returning the quantile boundaries in the distribution \mathcal{D}_i , split in g fragments. Also, q_l and q_u corresponds to the lower and upper quantile bounds, respectively. Then, the embedding of the mark m is performed according to $\mathcal{E}(m, v) = v'$, being \mathcal{E} the embedding function given in [11], and v' the resulting distorted value. Finally, if $v' \notin [q_l + k, q_u - k]$, embedding is rolled back and the algorithm proceeds checking the rest of R .

The main downside of this action is the WM capacity reduction (if WM length is too high with respect to η) due to rolling back the embedding of some marks. Nevertheless, WM recognition will be carried out as long as the number of tuples in R is higher than WM length.

3.2 Second Action: Change the Target Distribution

The second action is focused on saving those marks rolled back from the embedding in the previously described action (cf. Sect. 3.1). The attributes in the

selected tuple will be presented as a cyclic structure where A_{v-1} will precede A_0 (being A_i the i^{th} attribute of R). Then, if the value v' for A_i is out of its quantile, the embedding is rolled back, and the attribute $A_{(i+1)}$ is selected for the embedding. Moreover, attribute values in the ranges $[q_l, q_l+k]$ and $[q_u-k, q_u]$ are not considered for the embedding since it is very likely that v' will belong to the same range.

Finally, values of ξ and k are selected according to $k \geq [\xi]_{10}$ (being $[\xi]_{10}$ the decimal notation of the number of *lsbs*). This way, high pseudo-random embedding (which increases the difficulty for attackers to compromise WM detection) and a significant distortion reduction (with respect to methods not fragmenting \mathcal{D}_i in quantiles) will be achieved. Precisely, capacity is maintained while distortion resulting from WM embedding is reduced both at the binary level of v' (by applying the strategy given in [11]) and at the statistical distribution level of each attribute used as carrier.

4 Experimental Results

In the following, we present the experimental evaluation of the quantile-based watermarking actions for distortion reduction formalized in Sect. 3. Moreover, we discuss their benefits and downsides.

4.1 Experimental Setup

The data set used to perform the embedding and extraction of the watermark was *Forest Cover Type* [5], consisting of 581,012 tuples with 54 numerical attributes. Each one of the actions discussed in Sect. 3 was implemented based on a client/server architecture. The client layer was developed with Java 1.8 programming language and Eclipse Integrated Development Environment (IDE) 4.20. For the server layer was used Oracle Database 18C engine with Oracle SQL Developer 20.4 as Database Management System (DBMS) IDE. The runtime environment was a 2.11 GHz Intel i5 PC with 16.0 GHz of RAM with Windows 10 Pro OS.

We compare our results with a technique developed by Pérez Gort et al. [11] based on the AHK algorithm [2] and Sardroudi & Ibrahim's approach [13]. As mentioned in Sect. 2, the watermarking technique discussed in [11] uses a binary image to generate the watermark being embedded into R , and extends marks embedding to multiple attributes per tuple without considering numerical distortion preservation.

Figure 3 depicts the watermark sources we used, which are the binary images of the Chinese character *Dǎo* (20×21 pixels) and of the character *E* (10×10 pixels), respectively. Despite being binary images, missed pixels due to partial embedding, benign updates or attacks were highlighted using the red color for a clearer appreciation of the damage caused to the watermark.

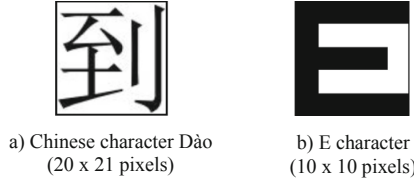


Fig. 3. Binary images used as watermark sources.

The metrics to analyze the quality of the extracted watermark with respect to the original image used for the WM generation were the correction factor (CF)¹ and the Structural Similarity Index (SSIM)² defined in [6].

$$M = |\mu - \mu'| \quad (2)$$

$$\Sigma = |\sigma - \sigma'| \quad (3)$$

The amount of distortion caused over each numerical attribute was measured by comparing the values of the mean μ and the standard deviation σ of the unwatermarked columns' numerical distribution with respect to the ones of the new distribution resulting from the embedding, denoted by μ' and σ' , respectively. Note that WM embedding allowing absolute distribution preservation is achieved when $M = 0$ and $\Sigma = 0$.

Furthermore, for cases when two different distributions present similar values of μ and σ we used the Kullback-Leibler divergence (D_{KL}), as depicted in Eq. (4), where $P_{\mathcal{D}_i}$ and $P_{\mathcal{D}'_i}$ represent the discrete probability distributions of the columns \mathcal{D}_i and \mathcal{D}'_i respectively, and \mathcal{X} indicates the probability space on which the distributions are defined.

$$D_{KL}(P_{\mathcal{D}_i}||P_{\mathcal{D}'_i}) = \sum_{x \in \mathcal{X}} P_{\mathcal{D}_i}(x) \log \left(\frac{P_{\mathcal{D}_i}(x)}{P_{\mathcal{D}'_i}(x)} \right) \quad (4)$$

4.2 Watermark Capacity Variations

The first requirement analyzed, featuring WM, is the capacity. A distortion reduction can be achieved by embedding fewer marks, which is not recommended since this will also reduce robustness.

Table 1 compares the capacity values obtained when the watermarking technique described in [11] is applied to the chosen data set, and when the same technique is enhanced by our actions. In particular, **NoQuant** captures the capacity when the quantile-based approach to watermark is not used, **NoEmb**

¹ CF $\in [0, 100]$ where 0 means total lack of correlation, and 100 the exact match between the extracted image with the original one.

² SSIM $\in [0, 1]$ where 0 represents the lack of similarity between the embedded and the extracted images, and 1 the presence of perfect similarity.

refers to the capacity obtained when mark embedding is canceled if quantile changes occur (cf. Sect. 3.1), and **Redist** to the capacity gained when distorted values are adjusted to prevent them from changing quantiles (cf. Sect. 3.2). For each case, the image of the synchronized WM and the correspondent SSIM and CF values are given.

Table 1. Watermark capacity varying γ .

γ	NoQuant		Proposals			
			NoEmb		Redist	
1						
	0.99	0.99	0.99	0.99	0.99	0.99
	99.76	99.00	99.76	99.00	99.76	99.00
5						
	0.99	0.99	0.99	0.99	0.99	0.99
	99.76	99.00	99.76	99.00	99.76	99.00
10						
	0.99	0.99	0.99	0.99	0.99	0.99
	99.76	99.00	99.76	99.00	99.76	99.00
20						
	0.99	0.99	0.97	0.99	0.99	0.99
	99.76	99.00	99.04	99.00	99.76	99.00
40						
	0.96	0.99	0.93	0.99	0.96	0.99
	97.14	99.00	95.71	99.00	97.14	99.00

The parameters' values for watermark synchronization were set as $SK = security2021$, $\delta = 5$, $\beta = 3$, and $\xi = 1$. Also, for the approaches fragmenting numerical distribution in quantiles we used $q = 4$ and $k = 1$. The experiments were carried out under a subset of *Forest Cover Type* data set composed by the first 30.000 tuples and 10 attributes.³

³ The subset selection was done to establish comparisons with other published results.

From the data reported in Table 1, can be concluded that: (i) The watermark capacity is not compromised when quantile-based embedding is performed (even for the line of action based on canceling mark insertion), and (ii) by using high γ values, the distortion caused by embedding is reduced without compromising the watermark recognition (especially for cases of watermark with small lengths).

In general, results achieved by canceling marks embedding experience a slight WM capacity reduction but, in terms of distortion, this strategy becomes highly recommended (especially when $\xi > 1$).⁴ Nevertheless, for preventing WM capacity reduction, the embedding of canceled marks in **NoEmb** is carried out in other locations on R by the strategy depicted in column **Redist**.

4.3 Imperceptibility Improvements

Regarding imperceptibility, by using $\xi = 1$ and $k = 1$ when applying the actions proposed in this work, there is evidence of a reduction of the distortion caused by the embedding in terms of M, preserving benefits and downsides in terms of Σ . Table 2 shows the values registered for M and Σ of each column of R, highlighting in blue color the results depicting lower distortion and in red color the ones causing more changes with respect to the approach not using quantiles.

Table 2. Distortion caused by WM embedding ($\gamma = 1, \xi = 1, k = 1$).

Attribute	NoQuant		Proposals			
	M	Σ	NoEmb		Redist	
			M	Σ	M	Σ
ATTR_01	8.50×10^{-3}	5.26×10^{-3}	8.20×10^{-3}	5.50×10^{-3}	8.23×10^{-3}	5.38×10^{-3}
ATTR_02	4.13×10^{-3}	5.01×10^{-4}	3.17×10^{-3}	1.08×10^{-3}	2.80×10^{-3}	6.33×10^{-4}
ATTR_03	7.00×10^{-3}	3.83×10^{-3}	1.19×10^{-2}	5.41×10^{-3}	2.90×10^{-3}	3.16×10^{-3}
ATTR_04	4.82×10^{-2}	1.25×10^{-2}	4.92×10^{-2}	1.46×10^{-2}	4.92×10^{-2}	1.46×10^{-2}
ATTR_05	4.93×10^{-3}	4.64×10^{-3}	5.50×10^{-3}	3.81×10^{-3}	4.53×10^{-3}	4.12×10^{-3}
ATTR_06	2.16×10^{-2}	9.56×10^{-3}	2.15×10^{-2}	9.57×10^{-3}	2.15×10^{-2}	9.62×10^{-3}
ATTR_07	2.56×10^{-2}	1.18×10^{-2}	1.64×10^{-2}	5.03×10^{-3}	1.84×10^{-2}	7.42×10^{-3}
ATTR_08	2.86×10^{-2}	1.42×10^{-2}	2.28×10^{-2}	8.11×10^{-3}	1.86×10^{-2}	9.31×10^{-3}
ATTR_09	8.30×10^{-3}	2.89×10^{-3}	7.60×10^{-3}	2.13×10^{-3}	6.93×10^{-3}	1.97×10^{-3}
ATTR_10	1.19×10^{-2}	2.34×10^{-3}	1.21×10^{-2}	2.50×10^{-3}	1.20×10^{-2}	2.53×10^{-3}

The presence of higher variation in some values of Table 2 is mainly due to the use of $\xi = 1$, which causes less distortion with respect to $k = 1$. Nevertheless, these values make the techniques vulnerable against bit flipping attacks,

⁴ The effect of the considered watermarking approaches over data distortion is discussed in Sect. 4.3.

being a perfect option for attackers to achieve WM removal without compromising data quality. Table 3 shows results by increasing the value of both ξ and k according to the recommendation given in Sect. 3. In this case, the robustness of our watermarking actions improves, whereas distortion experiments a significant reduction.

Table 3. Distortion caused by WM embedding ($\gamma = 1$, $\xi = 3$, $k = 4$).

Attribute	NoQuant		Proposals			
			NoEmb		Redist	
	M	Σ	M	Σ	M	Σ
ATTR_01	1.16×10^0	1.00×10^1	2.18×10^{-2}	9.96×10^{-3}	2.17×10^{-2}	1.00×10^{-2}
ATTR_02	6.85×10^{-1}	9.94×10^{-2}	2.77×10^{-3}	3.13×10^{-3}	4.33×10^{-3}	1.85×10^{-3}
ATTR_03	0	0	0	0	0	0
ATTR_04	6.16×10^{-2}	8.31×10^{-2}	6.23×10^{-3}	5.47×10^{-3}	6.23×10^{-3}	5.47×10^{-3}
ATTR_05	4.47×10^{-3}	9.87×10^{-3}	4.47×10^{-3}	9.87×10^{-3}	4.47×10^{-3}	9.87×10^{-3}
ATTR_06	3.08×10^0	4.32×10^0	2.30×10^{-2}	7.15×10^{-3}	2.28×10^{-2}	7.08×10^{-3}
ATTR_07	1.31×10^{-1}	1.37×10^0	1.45×10^{-2}	1.49×10^{-2}	1.30×10^{-2}	1.19×10^{-2}
ATTR_08	5.49×10^{-2}	1.21×10^0	1.17×10^{-2}	9.58×10^{-3}	1.10×10^{-2}	1.03×10^{-2}
ATTR_09	6.17×10^{-2}	2.56×10^{-1}	1.19×10^{-2}	7.52×10^{-3}	1.19×10^{-2}	7.25×10^{-3}
ATTR_10	6.57×10^{-1}	1.00×10^0	2.27×10^{-2}	7.29×10^{-3}	2.30×10^{-2}	7.57×10^{-3}

Table 4. Registered values of D_{KL} for experiments of Table 2.

Attribute	NoQuant	Proposals	
		NoEmb	Redist
ATTR_01	1.40×10^{-2}	1.39×10^{-2}	1.38×10^{-2}
ATTR_02	3.57×10^{-3}	3.54×10^{-3}	3.40×10^{-3}
ATTR_03	1.08×10^{-3}	1.78×10^{-3}	5.61×10^{-4}
ATTR_04	9.61×10^{-2}	9.09×10^{-2}	9.09×10^{-2}
ATTR_05	3.24×10^{-3}	3.29×10^{-3}	3.14×10^{-3}
ATTR_06	4.98×10^{-2}	4.98×10^{-2}	4.98×10^{-2}
ATTR_07	3.22×10^{-3}	1.99×10^{-3}	1.59×10^{-3}
ATTR_08	3.18×10^{-3}	2.80×10^{-3}	1.80×10^{-3}
ATTR_09	1.70×10^{-3}	1.88×10^{-3}	1.52×10^{-3}
ATTR_10	4.61×10^{-2}	4.60×10^{-2}	4.60×10^{-2}

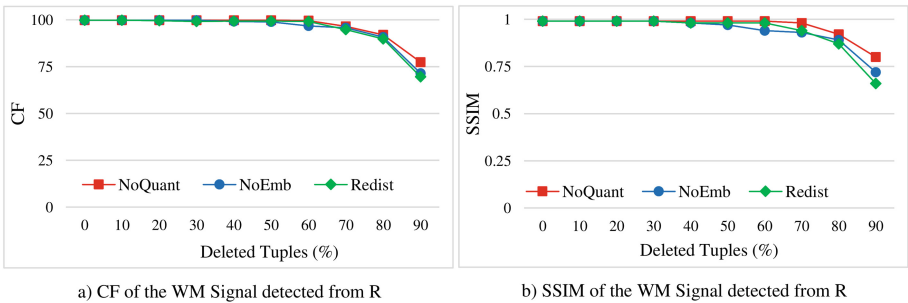
Table 5. Registered values of D_{KL} for experiments of Table 3.

Attribute	NoQuant	Proposals	
		NoEmb	Redist
ATTR_01	1.26×10^{-2}	1.19×10^{-2}	1.18×10^{-2}
ATTR_02	2.59×10^{-3}	2.37×10^{-3}	2.28×10^{-3}
ATTR_03	0	0	0
ATTR_04	6.83×10^{-2}	6.20×10^{-2}	6.20×10^{-2}
ATTR_05	1.67×10^{-3}	1.67×10^{-3}	1.67×10^{-3}
ATTR_06	4.91×10^{-2}	4.86×10^{-2}	4.86×10^{-2}
ATTR_07	2.53×10^{-3}	1.45×10^{-3}	9.87×10^{-4}
ATTR_08	2.42×10^{-3}	1.18×10^{-3}	7.08×10^{-4}
ATTR_09	2.19×10^{-3}	1.73×10^{-3}	1.28×10^{-3}
ATTR_10	4.42×10^{-2}	4.39×10^{-2}	4.39×10^{-2}

Tables 4 and 5 show the values of the D_{KL} metric for the experiments of Tables 2 and 3. The obtained results lead to the conclusion that the distributions resulting from applying the proposed lines of actions are more similar to the original data distributions than when the embedding is performed without considering quantiles.

4.4 Watermark Robustness Impact

Reducing distortion while preserving WM capacity has a positive impact on robustness. By performing the watermark embedding using $\gamma = 1$ and $\delta = 5$, all approaches guaranteed the WM signal total recovery for subset attacks based on inserting (or deleting) up to 90% of tuples with respect to the number of tuples stored in R. Instead, by using $\gamma = 10$, resilience against subset attacks will remain high. Nevertheless, because of WM capacity reduction, detected WM signal starts depicting small degradation when more than 80% of tuples are deleted (see Fig. 4).

**Fig. 4.** Quality of WM detected in R after different degree of subset deletion attacks.

Another feature of our strategies contributing to resilience against bit-flipping attacks is the increasing of the pseudo-random nature of WM embedding process. By selecting different numerical distributions in R, according to values in the database, and by increasing ξ and k , attackers face additional challenges for marks detection.

Besides the small variations in terms of robustness against subset deletion attacks for higher values of γ , a general appreciation in terms of WM capacity with respect to the distortion caused during WM embedding shows the benefits of proposed lines of action compared to traditional embedding. Figure 5 depicts the rate of WM quality (in terms of CF) vs. distortion. Considering that different attributes change values during the embedding, and that Fig. 5 reflects the whole distribution for each one of them, M_A and Σ_A were obtained from the average of M and Σ of all numeric columns used as carriers for each approach.

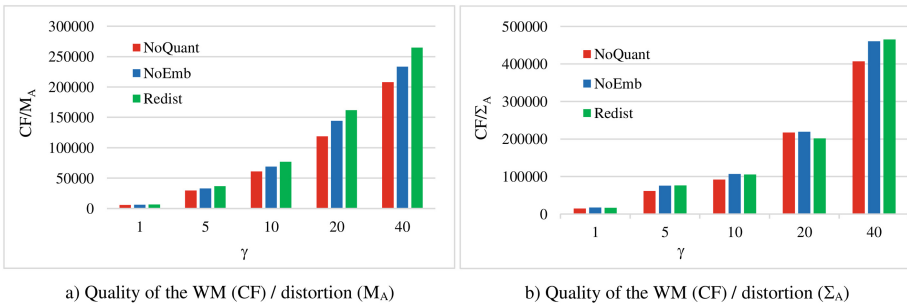


Fig. 5. Rate of detected WM quality/embedding distortion by varying γ .

4.5 Benefits of Selecting Meaningful Watermark Sources

Even for the action of rolling back mark embedding when quantile changes are spotted, WM capacity damages are not critical when WM length is not high, and meaningful WM sources are used. Table 6 shows the benefits obtained by considering symmetry criteria and neighboring pixels for the restoration of the extracted WM signal. Precisely, **PrevEnhancement** and **Enhancement** refers to the signal detected before and after the application of our enhancement actions, respectively. According to this behavior, by considering meaningful WM sources, rolling back mark embedding is another strategy worthy of being considered depending on the number of attributes and tuples being watermarked. In Table 6, the metric experimenting the increment regularly is the CF, which perceives the effects of recovering missed marks.

Table 6. WM signal enhancement (for meaningful WM sources).

γ	PrevEnhancement		Enhancement	
	SSIM	CF	SSIM	CF
20	0.97	99.04	0.98	99.28
40	0.93	95.71	0.91	97.14
60	0.81	82.85	0.83	91.66
80	0.78	79.76	0.74	90.23
100	0.59	64.04	0.59	84.52

5 Conclusions

In this paper, we proposed a quantile-based watermarking technique for relational data oriented to preserve the distribution of numerical attributes selected for mark embedding. Our approach follows two main lines of action: (i) rolling back mark embedding that violates quantile value preservation and (ii) seeking alternative embedding places for those marks causing a marked value changing quantile. Experimental results validate the relevance of *lsb* number and the threshold used for securing quantiles boundaries, for reducing the distortion while performing WM embedding. Furthermore, our technique shows an improvement in robustness while preserving WM capacity and increasing its imperceptibility.

Acknowledgement. This work has been partially supported by the project “VIR2EM - Virtualization and Remotization for Resilient and Efficient Manufacturing” - POR FESR VENETO 2014–2020.

References

1. Agrawal, R., Haas, P.J., Kiernan, J.: Watermarking relational data: framework, algorithms and analysis. *VLDB J.* **12**(2), 157–169 (2003)
2. Agrawal, R., Kiernan, J.: Watermarking relational databases. In: *VLDB 2002: Proceedings of the 28th International Conference on Very Large Databases*, pp. 155–166. Elsevier (2002)
3. Al-Haj, A., Odeh, A.: Robust and blind watermarking of relational database systems. *J. Comput. Sci.* **4**(12), 1024–1029 (2008)
4. Barni, M., Bartolini, F.: *Watermarking Systems Engineering: Enabling Digital Assets Security and Other Applications*. CRC Press, Boca Raton (2004)
5. Colorado-State-University: Forest CoverType, The UCI KDD Archive. Information and Computer Science. University of California, Irvine, June 1999. <http://kdd.ics.uci.edu/databases/covertyp/covertyp.html>
6. Gort, M.L.P., Olliaro, M., Cortesi, A., Uribe, C.F.: Semantic-driven watermarking of relational textual databases. *Expert Syst. Appl.* **167**, 114013 (2021)
7. Halder, R., Pal, S., Cortesi, A.: Watermarking techniques for relational databases: survey, classification and comparison. *J. Univers. Comput. Sci.* **16**(21), 3164–3190 (2010)

8. Hou, R., Xian, H.: A graded reversible watermarking scheme for relational data. *Mob. Netw. Appl.* 1–12 (2019)
9. Naz, F., et al.: Watermarking as a service (WaaS) with anonymity. *Multimedia Tools Appl.* **79**(23), 16051–16075 (2020)
10. Nematollahi, M.A., Vorakupipat, C., Rosales, H.G.: *Digital Watermarking: Techniques and Trends*. Springer, Heidelberg (2017)
11. Pérez Gort, M.L., Feregrino Uribe, C., Nummenmaa, J.: A minimum distortion: high capacity watermarking technique for relational data. In: *Proceedings of the 5th ACM Workshop on Information Hiding and Multimedia Security*, pp. 111–121 (2017)
12. Rani, S., Koshley, D.K., Halder, R.: Partitioning-insensitive watermarking approach for distributed relational databases. In: Hameurlain, A., Küng, J., Wagner, R., Dang, T.K., Thoai, N. (eds.) *Transactions on Large-Scale Data- and Knowledge-Centered Systems XXXVI*. LNCS, vol. 10720, pp. 172–192. Springer, Heidelberg (2017). https://doi.org/10.1007/978-3-662-56266-6_8
13. Sardroudi, H.M., Ibrahim, S.: A new approach for relational database watermarking using image. In: *5th International Conference on Computer Sciences and Convergence Information Technology*, pp. 606–610. IEEE (2010)
14. Siledar, S., Tamane, S.: A distortion-free watermarking approach for verifying integrity of relational databases. In: *2020 International Conference on Smart Innovations in Design, Environment, Management, Planning and Computing (ICSIDEMPC)*, pp. 192–195. IEEE (2020)
15. Sion, R., Atallah, M., Prabhakar, S.: Rights protection for relational data. *IEEE Trans. Knowl. Data Eng.* **16**(12), 1509–1525 (2004)
16. Sun, S., Xu, Y., Wu, Z.: Research on tampering detection of material gene data based on fragile watermarking. In: Sun, X., Wang, J., Bertino, E. (eds.) *ICAIS 2020*. CCIS, vol. 1252, pp. 219–231. Springer, Singapore (2020). https://doi.org/10.1007/978-981-15-8083-3_20
17. Xu, Y., Shi, B.: Copyright protection method of big data based on nash equilibrium and constraint optimization. *Peer-to-Peer Netw. Appl.* **14**(3), 1520–1530 (2021). <https://doi.org/10.1007/s12083-021-01096-4>
18. Zhao, M., Jiang, C., Duan, J.: Reversible database watermarking based on differential evolution algorithm. In: *2019 International Conference on Artificial Intelligence and Advanced Manufacturing (AIAM)*, pp. 120–124. IEEE (2019)

Open Access This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

