

UNIVERSIDAD DE LA REPÚBLICA

TESIS DE MAESTRÍA

---

# Coding of Multichannel Signals with Irregular Sampling

---

*Autor:*

*Pablo Cerveñansky*

*Supervisores:*

*Álvaro Martín*

*Gadiel Seroussi*

Núcleo de Teoría de la Información

Facultad de Ingeniería

September 29, 2020



# *Abstract*

TODO

# Table of Contents

<b>Abstract</b>	<b>ii</b>
<b>Table of Contents</b>	<b>iii</b>
<b>Index of Figures</b>	<b>iv</b>
<b>Index of Tables</b>	<b>vi</b>
<b>1 Datasets</b>	<b>2</b>
1.1 Overview . . . . .	3
1.2 Dataset IRKIS . . . . .	3
1.3 Dataset ADCP . . . . .	3
1.4 Dataset ElNino . . . . .	3
1.5 Dataset Solar . . . . .	3
1.6 Dataset Hail . . . . .	3
1.7 Dataset Tornado . . . . .	3
1.8 Dataset Wind . . . . .	3
<b>2 Algorithms</b>	<b>4</b>
2.1 Introduction . . . . .	5
2.2 Implementation details . . . . .	7
2.2.1 Gap Encoding in the Masking Variant . . . . .	9
2.3 Algorithm Base . . . . .	11
2.4 Algorithm PCA . . . . .	13
2.5 Algorithm APCA . . . . .	19
2.6 Algorithms PWLH and PWLHInt . . . . .	22
2.7 Algorithm SF . . . . .	30
2.8 Algorithm CA . . . . .	31
2.9 Algorithm FR . . . . .	37
2.10 Algorithm GAMPS . . . . .	42
<b>3 Experimental Results</b>	<b>43</b>
3.1 Experimental Setting . . . . .	44
3.2 Comparison of Masking and Non-Masking Variants . . . . .	46
3.3 Window Size Parameter . . . . .	50
3.4 Algorithms Performance . . . . .	54
3.4.1 Comparison with algorithm gzip . . . . .	58
3.5 Conclusions . . . . .	60
3.6 Future Work (TODO) . . . . .	60
<b>Bibliography</b>	<b>63</b>

# List of Figures

2.1	Coding pseudocode for the Constant and Linear model algorithms. . . . .	7
2.2	Decoding pseudocode for the Constant and Linear model algorithms. . . . .	8
2.3	Markov process diagram . . . . .	9
2.4	Base.code_column_ $NM$ pseudocode. . . . .	11
2.5	Base.decode_column_ $NM$ pseudocode. . . . .	11
2.6	PCA.code_column_ $M$ pseudocode. . . . .	13
2.7	PCA.decode_column_ $M$ pseudocode. . . . .	14
2.8	. . . . .	14
2.9	. . . . .	15
2.10	. . . . .	15
2.11	. . . . .	16
2.12	PCA.code_column_ $NM$ pseudocode. . . . .	17
2.13	PCA.decode_column_ $NM$ pseudocode. . . . .	18
2.14	APCA.code_column_ $M$ pseudocode. . . . .	19
2.15	APCA.decode_column_ $M$ pseudocode. . . . .	20
2.16	. . . . .	20
2.17	. . . . .	21
2.18	. . . . .	21
2.19	PWLH.code_column_ $M$ pseudocode. . . . .	22
2.20	PWLH.decode_column_ $M$ pseudocode. . . . .	23
2.21	. . . . .	24
2.22	. . . . .	24
2.23	. . . . .	25
2.24	. . . . .	25
2.25	. . . . .	26
2.26	. . . . .	26
2.27	. . . . .	27
2.28	. . . . .	27
2.29	. . . . .	28
2.30	. . . . .	28
2.31	Example SF . . . . .	30
2.32	CA.code_column_ $M$ pseudocode. . . . .	31
2.33	. . . . .	32
2.34	. . . . .	32
2.35	. . . . .	33
2.36	. . . . .	33
2.37	. . . . .	34
2.38	. . . . .	34
2.39	. . . . .	35
2.40	. . . . .	35
2.41	. . . . .	36
2.42	FR.code_column_ $M$ pseudocode. . . . .	37

2.43	FR.decode_column_ $M$ pseudocode. . . . .	38
2.44	push_points_indexes pseudocode. . . . .	39
2.45	. . . . .	39
2.46	. . . . .	40
2.47	. . . . .	40
2.48	. . . . .	41
3.1	CR and RD plots for every pair of algorithm variants $a_M, a_{NM} \in A_M$ , for the data type “SST” of the dataset SST. In the RD plot for algorithm PCA we highlight with a red circle the marker for the maximum value (50.60%) obtained for all the tested CAIs. . . . .	47
3.2	CR and RD plots for every pair of algorithm variants $a_M, a_{NM} \in A_M$ , for the data type “Longitude” of the dataset Tornado. In the RD plot for algorithm APCA we highlight with a blue circle the marker for the minimum value (-0.29%) obtained for all the tested CAIs. . . . .	48
3.3	Plots of $w_{global}^*$ , $w_{local}^*$ , and the RD between $c_{<a_v, w_{global}^*, e>}$ and $c_{<a_v, w_{local}^*, e>}$ , as a function of the threshold parameter $e$ , obtained for the data type “VWC” of the file “vwc_1202.dat.csv” of the dataset IRKIS. . . . .	51
3.4	Plots of $w_{global}^*$ , $w_{local}^*$ , and the RD between $c_{<a_v, w_{global}^*, e>}$ and $c_{<a_v, w_{local}^*, e>}$ , as a function of the threshold parameter $e$ , obtained for the data type “VWC” of the file “vwc_1203.dat.csv” of the dataset IRKIS. In the RD plot for algorithm PCA we highlight with a red circle the marker for the maximum value (10.68%) obtained for all the tested CAIs. . . . .	52
3.5	CR and window size parameter plots for every algorithm, for the data type “SST” of the dataset ElNino. For each threshold parameter $e \in E$ , we use blue circles to highlight the markers for the minimum CR value and the best window size parameter (in the respective plots corresponding to the best algorithm) . . . . .	55

# Índice de tablas

1.1	Datasets overview. The second column indicates the characteristic of each dataset, in terms of the amount of gaps. The third column shows the number of files. The fourth and fifth columns show the number of data types and their names, respectively. . . . .	3
2.1	Coding algorithms overview. For each algorithm we show whether it supports lossless and/or lossy compression (second and third columns), its model category (fourth column), whether the masking ( <i>M</i> ) and/or non-masking ( <i>NM</i> ) variants are valid (fifth and sixth columns), and whether the window size parameter ( <i>w</i> ) is supported (last column). . . . .	5
3.1	Range of values for the RD between the masking and non-masking variants of each algorithm (last column); we highlight the maximum (red) and minimum (blue) values taken by the RD. The results are aggregated by dataset. The second column indicates the characteristic of each dataset, in terms of the amount of gaps. The third column shows the number of cases in which the masking variant outperforms the non-masking variant of a coding algorithm, and its percentage among the total pairs of CAIs compared for a dataset. . . . .	49
3.2	RD between the OWS and LOWS variants of each CAI. The results are aggregated by algorithm and the range to which the RD belongs. . . . .	53
3.3	Compression performance of the best evaluated coding algorithm, for various error thresholds on each data type of each dataset. Each row contains information relative to certain data type. For each threshold, the first column shows the minimum CR, and the second column shows the base-2 logarithm of the window size parameter for the best algorithm (the one that achieves the minimum CR), which is identified by a certain cell color described in the legend above the table. . . . .	56
3.4	$\text{maxRD}(a, e)$ obtained for every pair of coding algorithm variant $a_v \in V^*$ and threshold parameter $e \in E$ . For each $e$ , the cell corresponding to the $\text{minmaxRD}(a)$ value is highlighted. . . . .	57
3.5	Compression performance of the best evaluated coding algorithm, for various error thresholds on each data type of each dataset, including the results obtained by gzip. Each row contains information relative to certain data type. For each threshold, the first column shows the minimum CR, and the second column shows the base-2 logarithm of the window size parameter for the best algorithm (the one that achieves the minimum CR), which is identified by a certain cell color described in the legend above the table. Algorithm gzip doesn't have a window size parameter, so the cell is left blank in these cases. . . . .	59
3.6	$\text{maxRD}(a, e)$ obtained for every pair of coding algorithm variant $a_v \in V^* \cup \{\text{gzip}\}$ and threshold parameter $e \in E$ . For each $e$ , the cell corresponding to the $\text{minmaxRD}(a)$ value is highlighted. . . . .	60

**Cambios de la versión anterior (13/7/2020) a esta versión (9/9/2020)****Chapter 1: Datasets**

- Agregué una introducción especificando lo que va a ir en cada section.
- En la Section 1.1 agregué una tabla con los datasets. Incluye la columna Dataset Characteristic, que también se muestra en la Tabla 3.1 del Chapter 3. Creo que está bueno que esa misma información haya sido mostrada antes en el informe (y de la misma manera) cuando se presentan los datasets.

**Chapter 2: Algorithms**

- Agregué una introducción especificando lo que va a ir en cada section.
- En la Section 2.1 agregué una tabla con los algoritmos. Creo que ayuda a entender el tema de los distintos conjuntos de algoritmos y variantes que se definen en el Chapter 3.
- En Section 2.2 agregué el pseudocódigo genérico para todos los algoritmos de modelo constante y lineal. El algoritmo GAMPS también considera la correlación entre columnas, así que el pseudocódigo es diferente a los demás (todavía no lo hice, lo pensaba poner en la section de GAMPS).
- Además de las primeras dos secciones, las secciones de Base y PCA (2.3 y 2.4) también están prontas para leer.
- En las secciones de APCA, CA y PWLH (2.5, 2.8 y 2.6) agregué el pseudocódigo del codificador con máscara y las imágenes del paso a paso del ejemplo. Todavía me falta escribir. No estoy seguro de que sea necesario agregar los pseudocódigos del decodificador con máscara. Tampoco sé si es necesario agregar los pseudocódigos del codificador y decodificador sin máscara en cada caso (ya lo hice con el algoritmo PCA y las diferencias son similares), a lo mejor alcanza con escribir un comentario en cada caso y listo.
- En las secciones de SF y FR (2.7 y 2.9) solamente agregué la imagen con la gráfica del último paso del ejemplo (Notar que para todos los algoritmos del capítulo, el ejemplo de 12 valores devuelve algo distinto al ser decodificado). Me falta todo el resto.
- En la sección de GAMPS (2.10) es la única en donde todavía no he escrito nada. Tengo que pensar otro ejemplo ya que este algoritmo codifica más de una columna del csv en paralelo.

**Chapter 3: Experimental Results**

- Introducción: Hice las correcciones marcadas y la separé en párrafos. Marqué un par de dudas con negrita.
- Al principio de la Section 3.1 agregué varias definiciones en un párrafo. No estoy seguro si debería hacer definiciones formales o si está bien que estén todas juntas en un párrafo.
- En las leyendas de las figuras 3.3 y 3.4 está bien poner LOWS y OWS, o debería poner  $w_{global}^*$  y  $w_{local}^*$ ?
- En todas las figuras después de la Section 3.2 creo que debería poner  $PCA_M$ ,  $APCA_M$ , etc. en vez de PCA, APCA, etc.
- Agregué la Section 3.5 con conclusiones. Algunas cosas quedaron similares a la introducción, pero intenté de que nunca hubiera dos frases iguales.
- Agregué la Section 3.6 con un punteo de ideas de trabajo futuro. Conclusions and Future Work debería ser un capítulo aparte?



# Chapter 1

## Datasets

In this chapter we present all of the datasets that were compressed in our experimental work, which is described in Chapter 3. In Section 1.1 we give an overview of the different datasets, describing their characteristics in terms of the amount of gaps, and the number of files and data types that each have. In the remaining sections we present in detail each of the datasets.

## 1.1 Overview

TODO:

- Explain why / how every dataset was transformed into a common format.
- Show example csv (describe header, column names, data rows, etc.)

Dataset	Dataset Characterstic	#Files	#Types	Data Types
IRKIS [1]	Many gaps	7	1	VWC
SST [2]	Many gaps	3	1	SST
ADCP [2]	Many gaps	3	1	Vel
Solar [3]	Many gaps	4	3	GHI, DNI, DHI
ElNino [4]	Many gaps	1	7	Lat, Long, Zonal Winds, Merid. Winds, Humidity, Air Temp., SST
Hail [5]	No gaps	1	3	Lat, Long, Size
Tornado [5]	No gaps	1	2	Lat, Long
Wind [5]	No gaps	1	3	Lat, Long, Speed

TABLE 1.1: Datasets overview. The second column indicates the characteristic of each dataset, in terms of the amount of gaps. The third column shows the number of files. The fourth and fifth columns show the number of data types and their names, respectively.

## 1.2 Dataset IRKIS

## 1.3 Dataset ADCP

## 1.4 Dataset ElNino

## 1.5 Dataset Solar

## 1.6 Dataset Hail

## 1.7 Dataset Tornado

## 1.8 Dataset Wind

## Chapter 2

# Algorithms

In this chapter we present all of the coding algorithms implemented in the project. In Section 2.1 we give an overview of the different algorithms, describing their features, parameters and masking variants. In Section 2.2 we provide some implementation details, including the pseudocodes for the coding and decoding subroutines that are common for every algorithm. We also describe the implementation of the gap encoding in the masking variant, which applies the KT estimator and arithmetic coding. In Section 2.3 we present algorithm Base, which is a trivial algorithm that serves as a base ground for the compression performance comparison developed in Chapter 3. In the remaining sections we present each of the coding algorithms in detail, including implementation specifics with the coding and decoding pseudocode for each of their variants, and an example that shows the encoding process step by step.

## 2.1 Introduction

There exists literature analyzing the performance of the state-of-the-art algorithms used for sensor data compression [6, 7]. In their original form, these algorithms assume that the signals have regular sampling and that there are no gaps in the data. However, it is often the case that this is not true for real-world datasets. For example, all of the datasets presented in Chapter 1 consist of signals that have either one or both of these characteristics. We selected a number of the state-of-the-art algorithms and implemented them in a way such that they are also able to encode signals with these characteristics.

The state-of-the-art algorithms follow a model-based compression approach: they represent the signal data using conventional approximation models whose parameters are then stored as compressed data. These model-based techniques manage to compress the data by exploiting their temporal and (in some cases) spatial correlation. They are popular not only due to their efficient compression performance, but also because they offer many benefits to data processing. For example, they can be used for inferring uncertain sensor readings, detecting outliers, indexing, etc. [6]. The model-based techniques are classified into four different categories, depending on their model type: *Constant models* encode signals using constant functions, *Linear models* use linear functions, *Nonlinear models* use complex nonlinear functions, and *Correlation models* simultaneously encode multiple signals while reducing redundant information. Algorithms in the last category are the only ones that, besides exploiting temporal correlation in the data, also exploit its spatial correlation.

In Table 2.1 we outline the features of every one of the algorithms that were implemented in the project. For each algorithm, the second and third columns indicate whether they support lossless and/or lossy compression, the fourth column shows its model category, the fifth and sixth columns indicate if the masking ( $M$ ) and/or non-masking ( $NM$ ) variants are valid, and the last column specifies if the window size parameter ( $w$ ) is supported.

Algorithm	Lossless	Lossy	Model	$M$	$NM$	$w$
Base	x		Constant		x	
PCA [8]	x	x	Constant	x	x	x
APCA [9]	x	x	Constant	x	x	x
PWLH [10]/ PWLHInt	x	x	Linear	x	x	x
SF [11]	x	x	Linear	x		
CA [12]	x	x	Linear	x	x	x
FR [13]	x	x	Linear	x		x
GAMPS [14]	x	x	Correlation	x	x	x

TABLE 2.1: Coding algorithms overview. For each algorithm we show whether it supports lossless and/or lossy compression (second and third columns), its model category (fourth column), whether the masking ( $M$ ) and/or non-masking ( $NM$ ) variants are valid (fifth and sixth columns), and whether the window size parameter ( $w$ ) is supported (last column).

Algorithm Base is a trivial lossless algorithm that is used as a base ground for comparing the performance of the remaining algorithms, all of which support both lossless and lossy encoding while guaranteeing a maximum point-by-point error between the decoded and the original files. This maximum error threshold is specified to the coding routine via the  $\epsilon$  parameter. We implemented algorithms PWLHInt and GAMPSLimit by modifying the code of algorithms PWLH and GAMPS, respectively, and this is why they appear in the same row.

We implemented at least one algorithm for each of the four categories defined above, except for the Nonlinear model. We decided not to implement algorithms in this category since they do not guarantee a maximum point-by-point error, and they yield poor compression performances [6].

We implemented two variants, masking ( $M$ ) and non-masking ( $NM$ ), which differ in the way they handle the encoding of the gaps in the data. The  $M$  variant of an algorithm first encodes the position of all the gaps and then proceeds to encode the data values. On the other hand, the  $NM$  variant encodes the position of the gaps and the data values simultaneously. Implementation details are presented in the remaining sections of the current chapter. It should be pointed out that, for every algorithm, the gaps in the decoded file always must match the gaps in the original file, regardless of the value of the  $\epsilon$  parameter, which is only considered when encoding the data values. We didn't find an efficient way of implementing the  $NM$  variant for the algorithms SF and FR, so they only support the  $M$  variant. In Section 3.2 we compare the compression performance of both variants,  $M$  and  $NM$ , for every algorithm that supports them.

Most of the algorithms support a window size parameter, which we label  $w$ . This parameter is passed to the coding routine, and it establishes the size of the blocks in which the data are processed and encoded. In algorithm PCA, parameter  $w$  defines a *fixed block size*, while in the rest of the algorithms it defines the *maximum block size*. More details on how the data are processed in blocks can be found in the pseudocodes for the algorithms presented in this chapter.

## 2.2 Implementation details

All of the algorithms are implemented in C++. In every case, we followed the original designs specified in the respective publications and added the new features described in Section 2.1. Algorithms PWLH [10], SF [11] and GAMPS [14] apply certain complex mathematical functions, so when implementing these algorithms we decided to reuse part of the code from the framework linked in [6]. The remaining algorithms [8, 9, 12, 13] were implemented entirely on our own.

Figures 2.1 and 2.2 show the coding and decoding pseudocodes, respectively, for every one of the Constant and Linear model algorithms implemented in the project (recall this information from Table 2.1). Both pseudocodes have the same length, and lines with the same number perform opposite actions. Constant and Linear model algorithms only exploit the temporal correlation in the data, thus they iterate through the data columns and encode them independently. Since Correlation models also exploit the spatial correlation (i.e. the data columns are *not* encoded independently), the pseudocodes for algorithms GAMPS and GAMPSLimit are different to the ones shown in this section, and we present them in Section 2.10.

In Figure 2.1, the inputs for the coding routine are a csv data file with the same format as the datasets presented in Chapter 1, an integer value (*algo\_key*) that uniquely describes the selected algorithm, a key (*v*) that describes its variant (either *M* or *NM*), and the maximum error threshold ( $\epsilon$ ) and window size (*w*) parameters. The output is a binary file, which represents the input file encoded by the selected algorithm with the specified variant and parameters.

```

input : in: csv data file to be coded
         algo_key: integer value for the selected algorithm
         v: key for the variant
          $\epsilon$ : maximum error threshold
         w: window size
output: out: binary file encoded with the selected algorithm
1  out = new_binary_file()
2  out.code_base_2(algo_key, 8)
3  algo = get_algorithm(algo_key)
4  if algo.has_window_param? then
5  |   out.code_base_2(w - 1, 8)
6  end
7  out.code_header(in.header)
8  out.code_base_2(in.rows_count(), 24)
9  out.code_ts_column(in.ts_column)
10 if v == M then
11 |   out.code_gaps(in)
12 |   foreach column in in.data_columns do
13 | |   algo.code_column_M(column, out,  $\epsilon$ , w, in.ts_column)
14 |   end
15 else
16 |   foreach column in in.data_columns do
17 | |   algo.code_column_NM(column, out,  $\epsilon$ , w, in.ts_column)
18 |   end
19 end
20 out.close_file()

```

FIGURE 2.1: Coding pseudocode for the Constant and Linear model algorithms.

The aforementioned binary file and the variant key are the only inputs for the decoding routine, in Figure 2.2. This routine also requires parameters *algo\_key* and *w*, but they are implicit inputs since they are already encoded in the binary file (lines 2 and 5). The output is a csv data file with the same headers as the original file (line 7), where the maximum absolute difference between any pair of decoded and original data values is equal to  $\epsilon$ , and the positions of the data gaps are the same as in the original file.

```

input : in: coded binary file
        v: key for the variant
output: out: decoded csv data file
1  out = new_csv_file()
2  algo_key = in.decode_base_2(8)
3  algo = get_algorithm(algo_key)
4  if algo.has_window_param? then
5    | w = in.decode_base_2(8) + 1
6  end
7  out.decode_header(in)
8  rows_count = in.decode_base_2(24)
9  ts_column = out.decode_ts_column(in)
10 if v == M then
11   | out.decode_gaps(in)
12   | while not in.reached_eof? do
13     | algo.decode_column_M(in, out, w, ts_column)
14   | end
15 else
16   | while not in.reached_eof? do
17     | algo.decode_column_NM(in, out, w, ts_column)
18   | end
19 end
20 out.close_file()

```

FIGURE 2.2: Decoding pseudocode for the Constant and Linear model algorithms.

The timestamp column, which is comprised of integers, is the first column in every csv data file, and it's also the first column to be encoded (line 9). This is done using a lossless schema, in which every integer is coded independently and using a fixed number of bits. It is worth mentioning that this coding schema could be further optimized. However, our project is focused in studying the compression of the data columns (i.e. the rest of the columns in the data file), since these are the only ones that could potentially have gaps.

When  $v = M$ , the masking variant of the algorithm is executed. First, the positions of the gaps in all of the data columns are encoded (line 11); implementation details are shown in Subsection 2.2.1. Then, the `code_column_M` subroutine is executed for every data column (lines 12-14). This subroutine has a different implementation depending on the specific algorithm, and the respective pseudocodes, both for the coder and the decoder, are presented in the following sections. Notice that these subroutines only encode the integer values in the data columns, since the positions of the gaps were already encoded in the previous step.

On the other hand, when  $v = NM$ , the non-masking variant of the algorithm is executed. In this case, the `code_column_NM` subroutine is executed for every data column (lines 16-18). Again, this subroutine is different for each specific algorithm, and we present the respective pseudocodes, for the coder and the decoder, in the following sections. Note that these subroutines must encode both the integer values *and the position of the gaps* in the data columns.

Before finishing this section, it should be pointed out that some of the implementation details were omitted in the pseudocodes for clarity. For instance, the key for the variant is not actually passed as an argument to the coding and decoding executables. Instead, we used a C++ macro and compiled two different executable files, one for each variant. This approach is useful both for optimizing the code and for making it more readable. Also, the coding and decoding subroutines (for both variants) do not always require all the arguments, and which ones are expected varies depending on the specific algorithm. In particular, the timestamp column is only used by the Linear model algorithms, and whether arguments  $\epsilon$  and  $w$  are required depends on the algorithm features, which are outlined in Table 2.1. Besides, the actual implementation allows to set different maximum error thresholds to encode different data columns, so  $\epsilon$  is actually a vector with a maximum error threshold parameter for each column in the input csv.

### 2.2.1 Gap Encoding in the Masking Variant

As we recall from previous sections, the masking variant of an algorithm first encodes the position of all the gaps in the data using a lossless schema, and afterwards it encodes the data values, using a lossless or lossy schema based on the value of the  $\epsilon$  parameter. These two actions correspond to lines 11, and 12-14, respectively, of the coding and decoding pseudocodes presented in Figures 2.1 and 2.2.

We decided to encode the position of the data gaps using a general purpose lossless algorithm called Arithmetic Coding (AC) [15, 16]. This algorithm sequentially encodes the symbols produced by a source, and it becomes more efficient the closer its model distribution is to the empirical distribution of the source [17]. In our case, the source is binary, with “0” denoting the presence of an integer data value, and “1” denoting the presence of a gap. As we recall from Chapter 1, the positions of the gaps follow different patterns for different datasets, but in general the gaps occur in bursts, and the amount of gaps is considerably less than the amount of data values. With this in mind, we thought it fitting to model the source distribution through a first-order Markov process with the Krichevsky–Trofimov estimator [18], which we define next in its binary alphabet form.

**Definition 2.2.1.** Given a source  $\pi$  with alphabet  $A = \{0, 1\}$ , and a string  $s$  generated by  $\pi$ , with  $|s_0|$  zeroes and  $|s_1|$  ones, the *Krichevsky–Trofimov estimator (KT)* assigns the following estimates  $p_i(s)$  to the probability of each symbol  $i \in A$ :

$$p_0(s) = \frac{|s_0| + 1/2}{|s_0| + |s_1| + 1}, \quad p_1(s) = \frac{|s_1| + 1/2}{|s_0| + |s_1| + 1}. \quad (2.1)$$

The first-order Markov process has two states,  $S_0$  and  $S_1$ , where the current state is  $S_i$  iff the last read symbol is  $i \in A$ . We set  $S_1$  as the initial state. In Figure 2.3 we display the diagram for this Markov process. To calculate the probabilities defined by (2.1) we need to keep a pair of counters (one for counting zeroes and another one for counting ones) in each state.

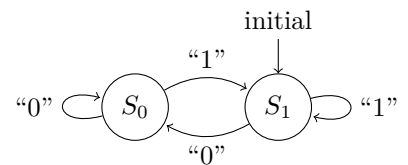


FIGURE 2.3: Markov process diagram

The version of the AC algorithm we used in our project is the CACM87 implementation [19, 20]. It is written in C and it is one of the most standard implementations. One of its advantages is that it allows to effortlessly set a custom model for the source. However, we had to overcome a minor obstacle to make it work within our scheme. In the CACM87 implementation, the coder closes the encoded file after it has encoded the last symbol. This implies that the decoder



recognizes that there are no more symbols left to decode once it reads the last byte of the encoded file. But this is not the case in our masking variant scheme, since after the AC coder has encoded the position of all the gaps in the data, our coding algorithm still has to encode all the data values before closing the encoded file. The problem materialized in the decoding process, because after the AC decoder had decoded the last byte corresponding to the position of the gaps (i.e. the last byte encoded by the AC coder), our decoding algorithm would occasionally continue processing bytes corresponding to the encoded data values, which naturally resulted in an error. The solution we found was to flush the current byte in the stream, before and after executing the AC algorithm, both in the coding and the decoding routines.

It is worth mentioning that we also experimented using Golomb coding [21] to encode the position of the gaps in the masking variant. We did this for every dataset and compared the results with those obtained by the AC algorithm. In most of the cases, the AC algorithm outperforms Golomb coding, and the few cases in which Golomb coding obtains a better compression, the difference is relatively small. Therefore, we decided to make AC the default algorithm for encoding the gaps in the masking variant.

## 2.3 Algorithm Base

Algorithm Base is a trivial lossless algorithm that serves as a base ground for comparing the performance of the rest of the algorithms. In particular, we reference it when defining the compression ratio metric (see Definition 3.1.4), which we use to assess the compression performance of an algorithm in Chapter 3.

In Figure 2.4 we show the pseudocode for the `code_column_M` subroutine. This subroutine iterates through every entry in the column of the csv data file, which is one of the inputs. Since algorithm Base only supports the *NM* variant, these entries can be either the character “N”, which represents a gap in the data (line 3), or an integer value representing an actual data measurement (line 5). Every column entry is encoded independently and using a fixed number of bits, namely `column.total_bits` (line 7), which varies *for each data type of each dataset*. In every case, the number of bits used for encoding a data value ultimately depends on the range and accuracy of the instrument that was used for measuring and storing the data. Therefore, without loss of generality, we can assume that `column.total_bits` is known by both the coding and decoding subroutines. The same can be said regarding `column.no_data_int`, a special integer reserved for encoding a gap, and `column.offset`, an offset that is added to the original value of a column entry to obtain a non-negative integer that can be encoded in base 2.

```

input : column: column of the csv data file to be coded
          out: binary file encoded with algorithm Base
1 foreach entry in column.entries do
2   if entry == “N” then
3     value = column.no_data_int
4   else
5     value = entry + column.offset
6   end
7   out.code_base_2(value, column.total_bits)
8 end

```

FIGURE 2.4: Base.`code_column_NM` pseudocode.

The pseudocode for the `decode_column_NM` subroutine is shown in Figure 2.5. This subroutine keeps running until every entry in the column has been decoded (i.e. the condition in line 1 becomes false). For decoding an entry, the first step is reading a fixed number of bits from the input binary file (line 2). Depending on the value that is read, either the character “N” (line 4) or an offset integer (line 6) is written into the decoded csv data file.

```

input : in: binary file coded with algorithm Base
          out: decoded csv data file
1 while not in.column_decoded? do
2   value = in.decode_base_2(column.total_bits)
3   if value == column.no_data_int then
4     out.write_string(“N”)
5   else
6     out.write_string(value - column.offset)
7   end
8 end

```

FIGURE 2.5: Base.`decode_column_NM` pseudocode.

The details regarding the calculation of the conditional in line 1 of the decoding subroutine were omitted in the pseudocode for clarity. However, the total number of entries in a column is known.

by the decoder (recall line 8 in the decoding subroutine in Figure 2.2), so all it has to do is keep a counter, update it every time an entry is decoded, and check its value in every iteration. This explanation is also valid for the remaining decoding subroutines presented in this chapter.

It is worth mentioning that, in the following sections, with the purpose of simplifying the pseudocodes and the descriptions of the algorithms, we always omit the offset operations. Yet the pertinent comments made in this section apply to every algorithm.

## 2.4 Algorithm PCA

Algorithm PCA [8], also known as Piecewise Constant Approximation, supports lossless and lossy compression, with both variants ( $M$  and  $NM$ ), and it has a window size parameter ( $w$ ) that establishes a fixed block size in which the data are processed and encoded. It is a Constant model algorithm, so it encodes signals using constant functions.

In Figure 2.6 we show the pseudocode for the `code_column_M` subroutine. Since this subroutine corresponds to the masking variant, we can assume that all of the processed column entries are integer values. We observe that these column entries are added to the window (line 3) until its size is equal to  $w$  (i.e. the condition in line 4 becomes false). When that happens, the window is encoded, which can occur in two different ways. If the absolute difference between the maximum and minimum values is less than or equal to  $2 * \epsilon$  (i.e. the condition in line 7 is true), then bit 0 and the mid-range of the window are encoded (lines 8-10). On the other hand, if the condition in line 7 is false, then bit 1 and each of the window values are encoded (lines 12-15). After a window is encoded, a new one is created (line 17), and the process is repeated until every entry in the column of the csv data file has been processed.

```

input : column: column of the csv data file to be coded
        out: binary file encoded with algorithm PCA
         $\epsilon$ : maximum error threshold
         $w$ : fixed window size
1  win = new_window()
2  foreach entry in column.entries do
3      win.push(entry)
4      if win.size <  $w$  then
5          continue
6      end
7      if |win.max - win.min| ≤  $2 * \epsilon$  then
8          out.code_bit(0)
9          mid_range = (win.min + win.max)/2
10         out.code_base_2(mid_range, column.total_bits)
11      else
12          out.code_bit(1)
13          foreach win_val in win.values do
14              out.code_base_2(win_val, column.total_bits)
15          end
16      end
17      win = new_window()
18 end

```

FIGURE 2.6: PCA.code\_column\_M pseudocode.

It should be pointed out that, since the number of column entries is not always divisible by  $w$ , it is possible that in the last iteration the last window is not encoded. Therefore, after the last iteration the algorithm must check whether the window is empty, and if that is not the case it must encode it. In this scenario, the last window is encoded by executing the same code as in lines 12-15, which was left out of the pseudocode for clarity.

The pseudocode for the `decode_column_M` subroutine is shown in Figure 2.7. This subroutine keeps running until every entry in the column has been decoded. First, a single bit is read from the input binary file (line 2). If that bit is 0, then a fixed number of bits is read next, and the associated base-2 integer, which in the coding process corresponded to the mid-range of a window, is written into the decoded csv data file  $w$  times (lines 4-7). On the other hand, if the

bit read is 1, then the following process is repeated a total of  $w$  times: a fixed number of bits is read, and the associated base-2 integer is written into the decoded csv data file (lines 9-12).

```

input : in: binary file coded with algorithm PCA
         out: decoded csv data file
         w: fixed window size
1 while not in.column_decoded? do
2   bit = in.decode_bit()
3   if bit == 0 then
4     mid_range = in.decode_base_2(column.total_bits)
5     repeat w times
6       out.write_string(mid_range)
7     end
8   else
9     repeat w times
10    value = in.decode_base_2(column.total_bits)
11    out.write_string(value)
12  end
13 end
14 end

```

FIGURE 2.7: PCA.decode\_column\_ $M$  pseudocode.

Next we present an example that illustrates the mechanics of algorithm PCA. In particular, we describe the process of encoding a data column that consists of 12 data values, where each data value was recorded at a certain timestamp. The data value vs. timestamp graph is shown in Figure 2.8. The blue crosses represent the original data values. As the encoding process is completed, the red circles will represent the encoded data values, which are the values written into the coded binary file by the code\_column\_ $M$  subroutine, and the orange dots will represent the decoded data values, which are the values written into the decoded csv file by the decode\_column\_ $M$  subroutine.

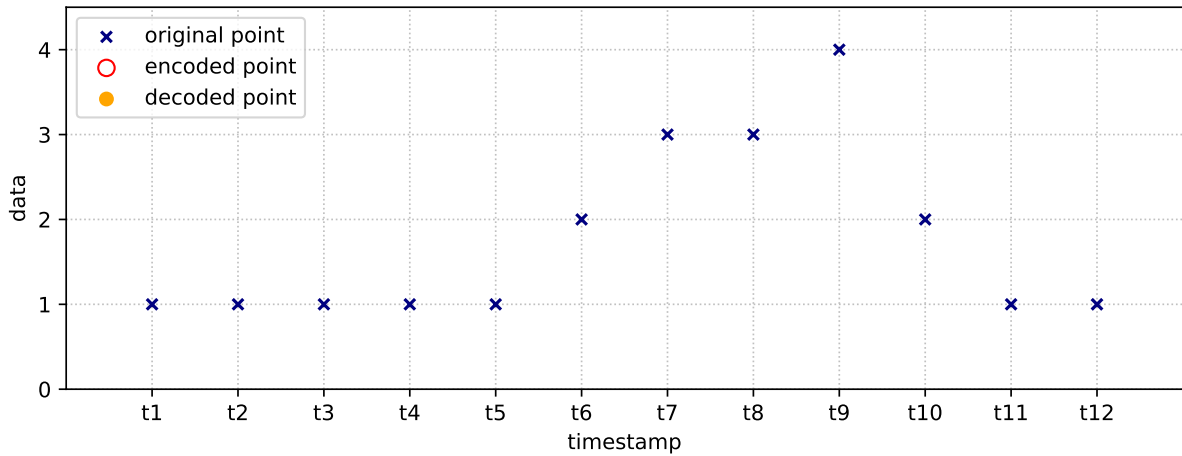


FIGURE 2.8

We consider that the difference between any pair of adjacent timestamps is equal to 60. This information is relevant, since we use the same data column in the following sections to illustrate the mechanics of the remaining algorithms, including the Linear model algorithms, which must acknowledge the timestamp values in the encoding process. However, since PCA is a Constant

model algorithm, the specific timestamp values are irrelevant in the current example. In this example we use algorithm PCA with an error threshold parameter ( $\epsilon$ ) equal to 1, and a fixed window size ( $w$ ) equal to 4.

Since  $w = 4$  and the first four data values in the column are all equal (to 1), the condition in line 7 of the `code_column_M` subroutine is true, because  $|1 - 1| \leq 2 * 1$ . Therefore, the first window is encoded by executing lines 8-10. In this case, the mid-range of the window is 1, so the first four data values are encoded as 1. Figure 2.9 shows this step in the graph. Notice that, since all of the values in the window are equal, in this case the condition in line 7 would be true regardless of the value of parameter  $\epsilon$ .

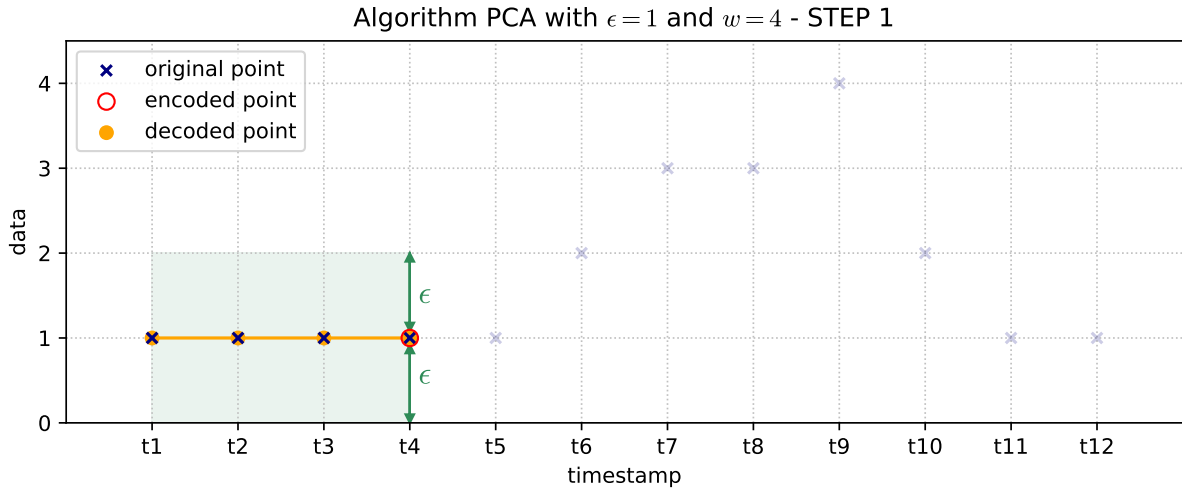


FIGURE 2.9

The second window is filled with the next four data values in the column, i.e.  $[1, 2, 3, 3]$ . Again, the condition in line 7 is true, since  $|3 - 1| \leq 2 * 1$ , but in this case the mid-range is 2, so these four data values are encoded as 2. This step is shown in Figure 2.10.

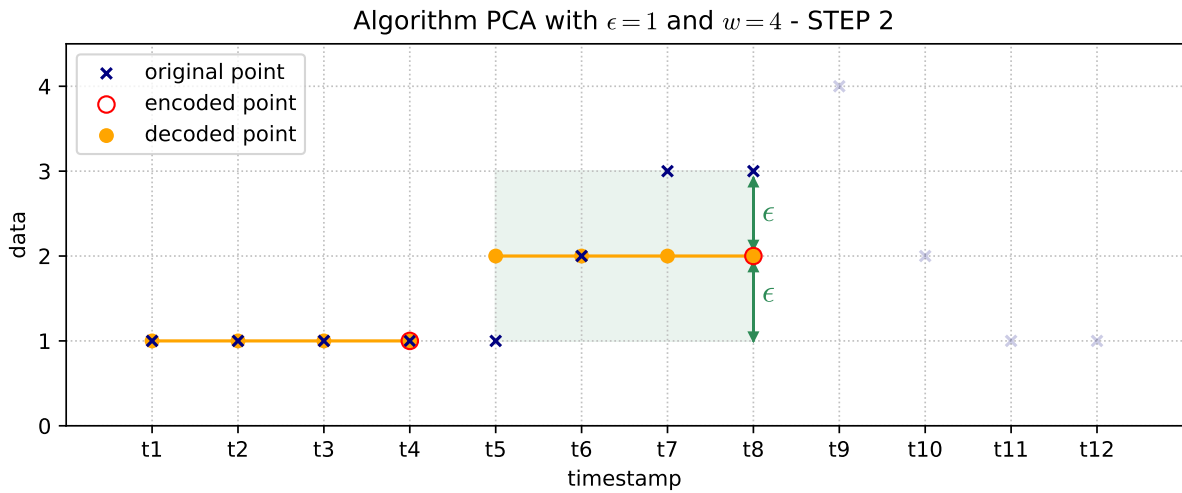


FIGURE 2.10

The third and last window is filled with the last four data values in the column, i.e.  $[4, 2, 1, 1]$ . In this case, the condition in line 7 is false, since  $|4 - 1| > 2 * 1$ , so the window is encoded by executing lines 12-15, which means that each of the data values in the window is encoded independently with its original value. This last step is shown in Figure 2.11.

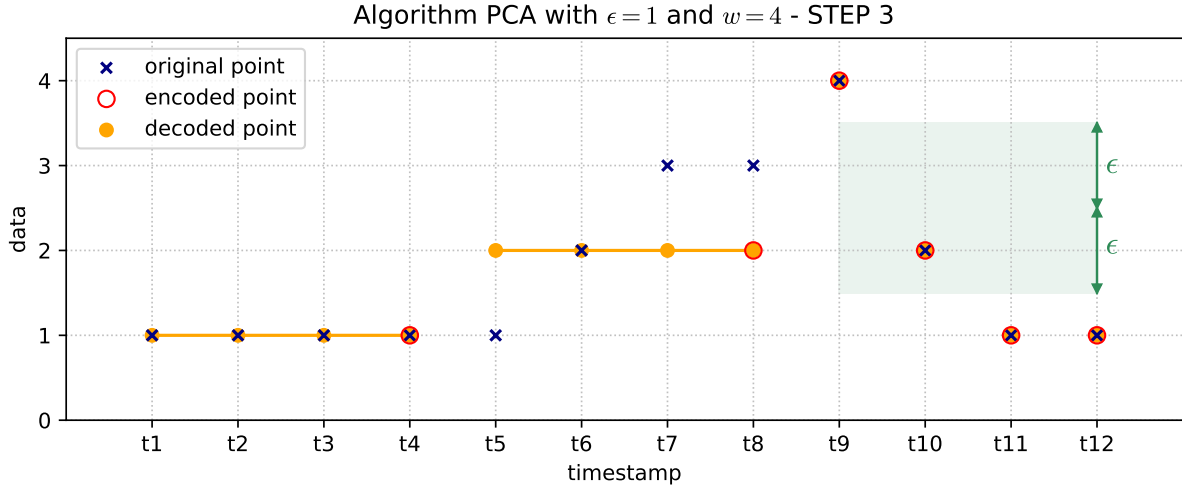


FIGURE 2.11

This simple example fairly represents every scenario that might manifest in the encoding process. Since the threshold condition holds for the first two windows, both are encoded with exactly the same amount of bits, i.e.  $1 + \text{column.total\_bits}$ , with the actual bits differing due to the fact that the mid-range is not the same in both cases. On the other hand, since the threshold condition does not hold for the last window, it is encoded with  $1 + w * \text{column.total\_bits}$  bits. This example serves to show why algorithm PCA is expected to achieve better compression performances on constant rather than variable data sources.

So far in this section, we have focused on the masking variant of algorithm PCA. In what follows, we detail the features of its non-masking variant. In Figure 2.12 we show the pseudocode for the `code_column_NM` subroutine. In this case, the processed column entries can be, not only an integer value representing an actual data measurement, but also the character “N”, which represents a gap in the data. The first six lines are the same as in the `code_column_M` subroutine: the column entries are added to the window until its size is equal to  $w$ . When that happens, the window is encoded, which can occur in three different ways. If all of the window entries are integers and the absolute difference between the extremes is less than or equal to  $2 * \epsilon$  (i.e. the condition in line 7 is true), then bit 0 and the mid-range of the window are encoded (lines 8-10). If instead all of the window entries represent a gap in the data (i.e. the condition in line 11 is true), then bit 0 and `column.no_data_int`, a special integer reserved for encoding a gap, are encoded (lines 12-13). Observe that in these first two cases, the windows are encoded with exactly the same amount of bits, i.e.  $1 + \text{column.total\_bits}$ . If none of the previous two conditionals are true for the window, then bit 1 and each of the window entries are encoded, again using `column.no_data_int` for encoding the gaps (lines 15-19). In this last case, the window is encoded with  $1 + w * \text{column.total\_bits}$  bits.

```

input : column: column of the csv data file to be coded
        out: binary file encoded with algorithm PCA
         $\epsilon$ : maximum error threshold
        w: fixed window size
1  win = new_window()
2  foreach entry in column.entries do
3      win.push(entry)
4      if win.size < w then
5          continue
6      end
7      if win.all_entries_are_integers and  $|win.max - win.min| \leq 2 * \epsilon$  then
8          out.code_bit(0)
9          mid_range = (win.min + win.max)/2
10         out.code_base_2(mid_range, column.total_bits)
11     else if win.all_entries_are_no_data then
12         out.code_bit(0)
13         out.code_base_2(column.no_data_int, column.total_bits)
14     else
15         out.code_bit(1)
16         foreach win_val in win.values do
17             value = (win_val == "N") ? column.no_data_int : win_val
18             out.code_base_2(value, column.total_bits)
19         end
20     end
21     win = new_window()
22 end

```

FIGURE 2.12: PCA.code\_column\_NM pseudocode.

It should be clear that the pseudocodes for the code\_column\_NM and code\_column\_M subroutines are rather similar. The difference being that the code\_column\_NM subroutine must handle gaps in the data, which requires minor tweaks in the conditionals. The decode\_column\_NM subroutine, whose pseudocode is shown in Figure 2.13, is also quite similar to the decode\_column\_M subroutine. The only differences being the two additional lines (5 and 12) that are needed for decoding the gaps.



```

input : in: binary file coded with algorithm PCA
        out: decoded csv data file
        w: fixed window size
1 while not in.column_decoded? do
2   bit = in.decode_bit()
3   if bit == 0 then
4     value = in.decode_base_2(column.total_bits)
5     string = (value == column.no_data_int) ? "N" : value
6     repeat w times
7       | out.write_string(string)
8     end
9   else
10    repeat w times
11      | value = in.decode_base_2(column.total_bits)
12      | string = (value == column.no_data_int) ? "N" : value
13      | out.write_string(string)
14    end
15  end
16 end

```

FIGURE 2.13: PCA.decode\_column\_NM pseudocode.

## 2.5 Algorithm APCA

Algorithm APCA [9], also known as Adaptive Piecewise Constant Approximation, is a Constant model algorithm, so it encodes signals using constant functions. As its name suggests, it operates similarly to algorithm PCA, the difference being that in APCA the size of the blocks in which the data are processed and encoded is not fixed, but variable. The window size parameter ( $w$ ) establishes the maximum block size allowed for the algorithm. It supports lossless and lossy compression, with both variants ( $M$  and  $NM$ ).

In Figure 2.14 we present the pseudocode for the `code_column_M` subroutine. Since this subroutine corresponds to the masking variant, we assume that all of the processed column entries are integer values. In every iteration an entry is added to the window (line 3), and a conditional that depends on parameters  $w$  and  $\epsilon$  is checked (line 4). If the size of the window is less than or equal to  $w$ , and the absolute difference between the maximum and minimum window values is less than or equal to  $2 * \epsilon$ , then the algorithm continues with the next iteration. Otherwise, when the window violates any of the constraints, the entry is removed from the window (line 7), the window is encoded (lines 8-10), and a new window that includes the entry is created (lines 11-12). Observe that every window is encoded with the same amounts of bits, i.e.  $\log_2 w + \text{column.total\_bits}$ , where  $\log_2 w$  bits are used for encoding its size, and  $\text{column.total\_bits}$  bits are used for encoding its mid-range.

```

input : column: column of the csv data file to be coded
        out: binary file encoded with algorithm APCA
         $\epsilon$ : maximum error threshold
         $w$ : maximum window size
1 win = new_window()
2 foreach entry in column.entries do
3   win.push(entry)
4   if win.size  $\leq w$  and  $|win.max - win.min| \leq 2 * \epsilon$  then
5     continue
6   end
7   entry = win.unpush()
8   out.code_base_2(win.size - 1, log2 w)
9   mid_range = (win.min + win.max)/2
10  out.code_base_2(mid_range, column.total_bits)
11  win = new_window()
12  win.push(entry)
13 end

```

FIGURE 2.14: `APCA.code_column_M` pseudocode.

It should be noted that, after executing the last iteration, the window is never empty, which means that its values have not yet been encoded. Therefore, after the last iteration the algorithm must encode the window, which is done by executing the same code as in lines 8-10. This was left out of the pseudocode for clarity.

The pseudocode for the `decode_column_M` subroutine is shown in Figure 2.15. This subroutine is fairly simple. First, the window size (line 2) and its mid-range value (line 3) are decoded. Then, the mid-range value is written into the decoded csv data file times the window size (lines 4-6). The subroutine keeps running until every entry in the column has been decoded.

```

input : in: binary file coded with algorithm APCA
        out: decoded csv data file
        w: maximum window size
1 while not in.column_decoded? do
2   win_size = in.decode_base_2( $\log_2 w$ ) + 1
3   mid_range = in.decode_base_2(column.total_bits)
4   repeat win_size times
5     out.write_string(mid_range)
6   end
7 end

```

FIGURE 2.15: APCA.decode\_column\_M pseudocode.

Next we present an example that illustrates the mechanics of algorithm APCA. As we did with algorithm PCA, we describe the process of encoding the data values which are shown in Figure 2.8. The comments regarding the timestamp values and the graph legends that were made when introducing said figure, also apply in this case. In this example we use algorithm APCA with an error threshold parameter ( $\epsilon$ ) equal to 1, and a maximum window size ( $w$ ) equal to 256.

The condition in line 4 of the code `_column_M` subroutine is true for the first eight iterations, i.e.  $[1, 1, 1, 1, 2, 3, 3]$ , are added to the first window. The value processed in the 9th iteration is 4, which causes the second constraint to become false, since  $|4 - 1| > 2 * 1$ . Therefore, that value is removed from the first window and added to a new second window, and the first window is encoded, which requires  $\log_2 w = \log_2 256 = 8$  bits for encoding its size (i.e. 8), and *column.total\_bits* for encoding its mid-range (i.e. 2). This step is shown in Figure 2.16.

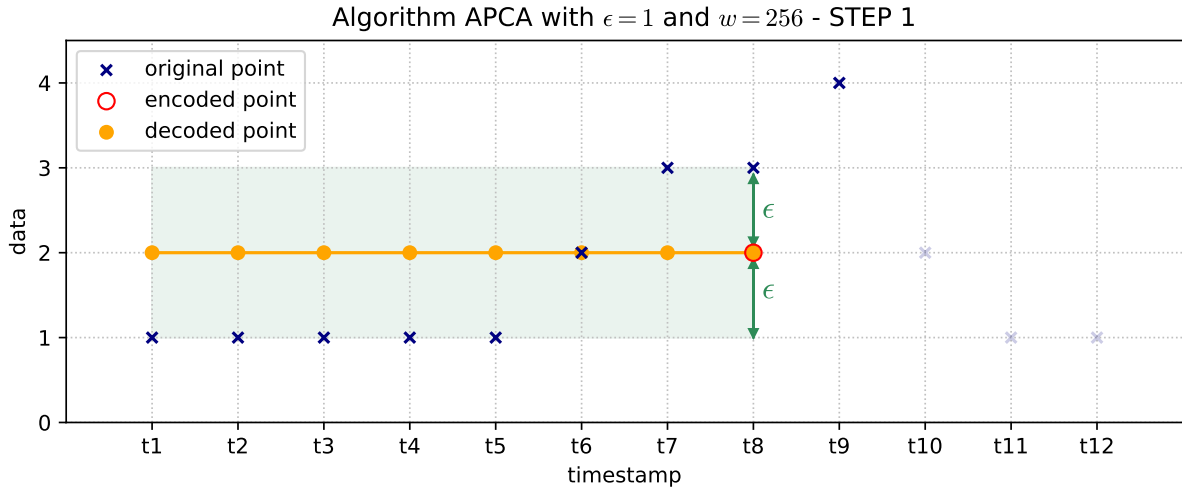


FIGURE 2.16

For the second window, the condition in line 4 is true in the 10th iteration. However, the second constraint becomes false in the 11th iteration, for data value 1, since again  $|4 - 1| > 2 * 1$ . That value is added to a new third window, and the second window is encoded. In this case the window is equal to  $[4, 2]$ , so its size is 2 and its mid-range is 3. This step is shown in Figure 2.17.

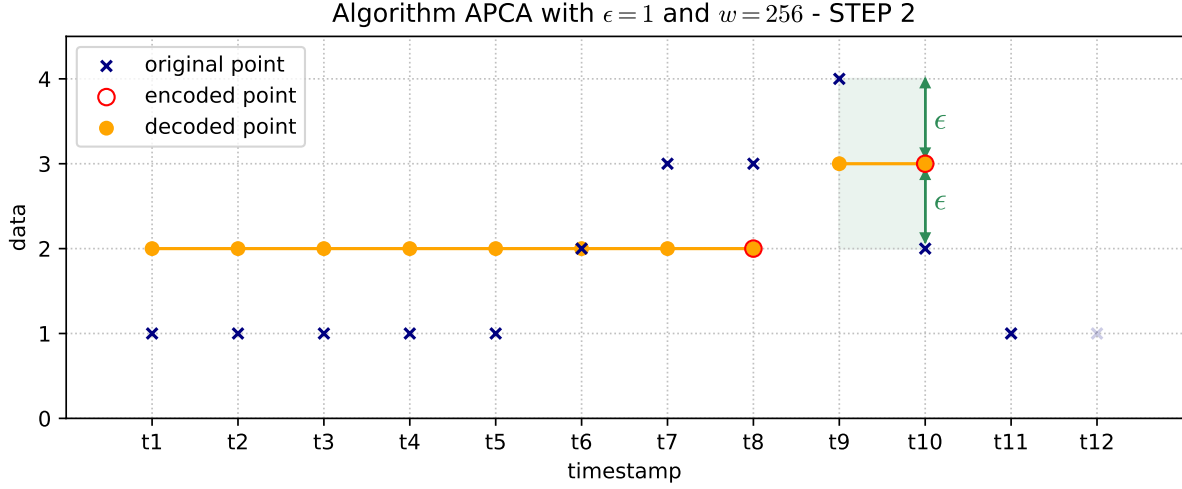


FIGURE 2.17

For the third window, the condition in line 4 is true for the 12th and last iteration. This means that said window, which is equal to  $[1, 1]$ , must be encoded after executing the last iteration. In this case, its size is 2 and its mid-range is 1. This last step is shown in Figure 2.18.

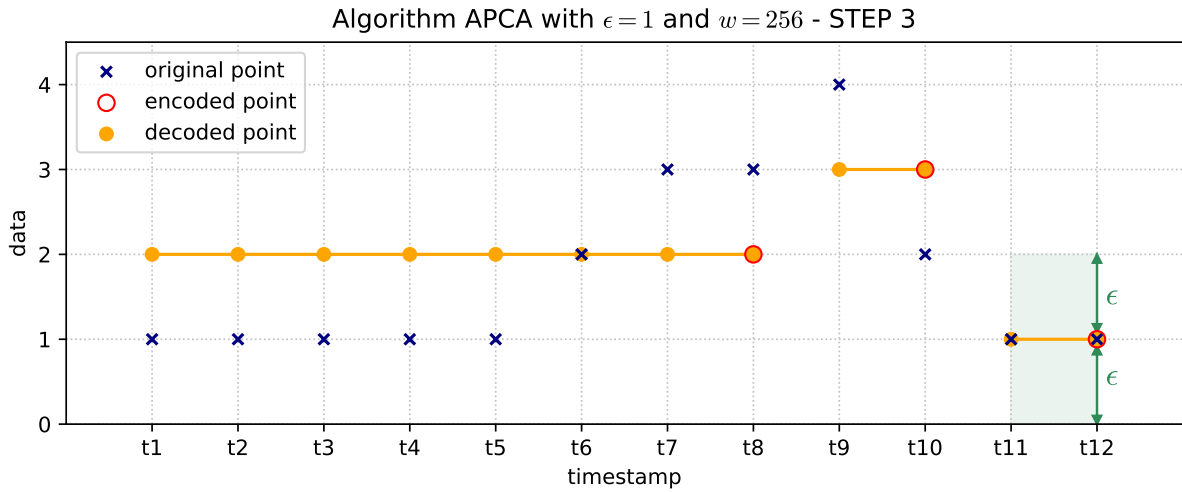


FIGURE 2.18

## 2.6 Algorithms PWLH and PWLHInt

Algorithm PWLH [10], also known as PieceWise Linear Histogram, supports lossless and lossy compression, with both variants ( $M$  and  $NM$ ), and it has a window size parameter ( $w$ ) that establishes the maximum block size in which the data are processed and encoded. It is a Linear model algorithm, so it encodes signals using linear functions. In particular, the data points in each window are modeled by a line segment that minimizes the maximum distance from these points to the segment. For the operations in the two-dimensional Euclidean space, which involve calculating said segment, and computing the convex hull of the data points by applying Graham's Scan algorithm [22], we reused the code from the framework linked in [6].

Figure 2.19 shows the pseudocode for the `code_column_M` subroutine for algorithm PWLH. Creating a new window always involves creating an associated empty convex hull (lines 1 and 24), and every time a data entry is added or removed from the window, the convex hull must be updated (lines 5, 13, 16 and 26). The window is coded in two scenarios: when it reaches the maximum size allowed (line 9), or when the convex hull violates the threshold condition (line 14). This threshold condition is valid iff there exists an edge in the convex hull for which the maximum distance from any of the points in the hull to said edge is less than or equal to  $2 * \epsilon$ .

```

input : column: column of the csv data file to be coded
        out: binary file encoded with algorithm PWLH
         $\epsilon$ : maximum error threshold
         $w$ : maximum window size
1  win = new_window()
2  foreach entry in column.entries do
3      if win.size == 0 then
4          win.push(entry)
5          win.update_convex_hull()
6          continue
7      end
8      code_window = false
9      if win.size ==  $w$  then
10         code_window = true
11     else
12         win.push(entry)
13         win.update_convex_hull()
14         if not win.PWLH_condition_holds?( $\epsilon$ ) then
15             entry = win.unpush()
16             win.update_convex_hull()
17             code_window = true
18         end
19     end
20     if code_window then
21         out.code_base_2(win.size - 1,  $\log_2 w$ )
22         point_A, point_B = win.get_approximation_segment()
23         out.code_float(point_A.y)
24         out.code_float(point_B.y)
25         win = new_window()
26         win.push(entry)
27         win.update_convex_hull()
28     end
29 end

```

FIGURE 2.19: PWLH.code\_column\_M pseudocode.

Encoding a window involves encoding its size (line 21) together with the y-coordinates of the two endpoints of the segment that minimizes the maximum distance from the points in the window to the segment (lines 23-24). The window size is encoded using  $\log_2 w$  bits, while the y-coordinates are encoded as float values, i.e. using 4 bytes, since this is the precision adopted in the method we reused for calculating said segment, which is invoked in line 22. Notice that the values of the x-coordinates were previously encoded with the timestamp column (recall line 9 in the pseudocode presented in Figure 2.1), so this coding subroutine must not encode them again. However, it is important to point out that, since PWLH is a Linear model algorithm, for the operations in the two-dimensional Euclidean space to make sense the x-coordinates must be considered in the actual coding routine, but they were omitted in the pseudocode for clarity. After encoding the window, a new window and its convex hull are created with the current entry (lines 25-27).

In Figure 2.20 we present the pseudocode for the `decode_column_M` subroutine. This subroutine keeps running until every entry in the column has been decoded. In each iteration a single window is decoded. First, the window size (line 2) and a pair of floats corresponding to the y-coordinates of the segment endpoints (lines 3-4) are decoded. Next, the approximation segment associated to the window is created (line 6), and the algorithm iterates through the index of every window entry, calculates its data value and writes it into the decoded csv data file (lines 7-9). The data value is obtained by replacing the x variable in the segment equation with the timestamp associated to the index of the window entry. Since the csv file must only consist of integer values, values are rounded to the nearest integer. Again, operations with the timestamps were omitted in the pseudocode for clarity, but they must also be considered in the decoding routine.

```

input : in: binary file coded with algorithm PWLH
        out: decoded csv data file
        w: maximum window size
1 while not in.column_decoded? do
2   win_size = in.decode_base_2( $\log_2 w$ ) + 1
3   point_A_y = in.decode_float()
4   point_B_y = in.decode_float()
5   win = new_window(win_size)
6   win.create_approximation_segment(point_A_y, point_B_y)
7   foreach index in [0..win_size - 1] do
8     value = win.segment_equation(index)
9     out.write_string(value)
10  end
11 end

```

FIGURE 2.20: PWLH.decode\_column\_M pseudocode.

Next we present an example that illustrates the mechanics of algorithm PWLH. Again, we describe the process of encoding the data values shown in Figure 2.8, where the distance between any pair of adjacent timestamps is equal to 60. In this example we use algorithm PWLH with an error threshold parameter ( $\epsilon$ ) equal to 1, and a maximum window size ( $w$ ) equal to 256.

Since there are only 12 data values to encode, no window in this example can reach the maximum size (256). Therefore, a window is only encoded when its convex hull violates the threshold condition in line 14. In the first iteration the window is empty, so the algorithm just adds the first data value to the window (lines 3-7). Figure 2.21 shows this step in the graph. Observe that, besides the original values, the convex hull for the current window is also shown.

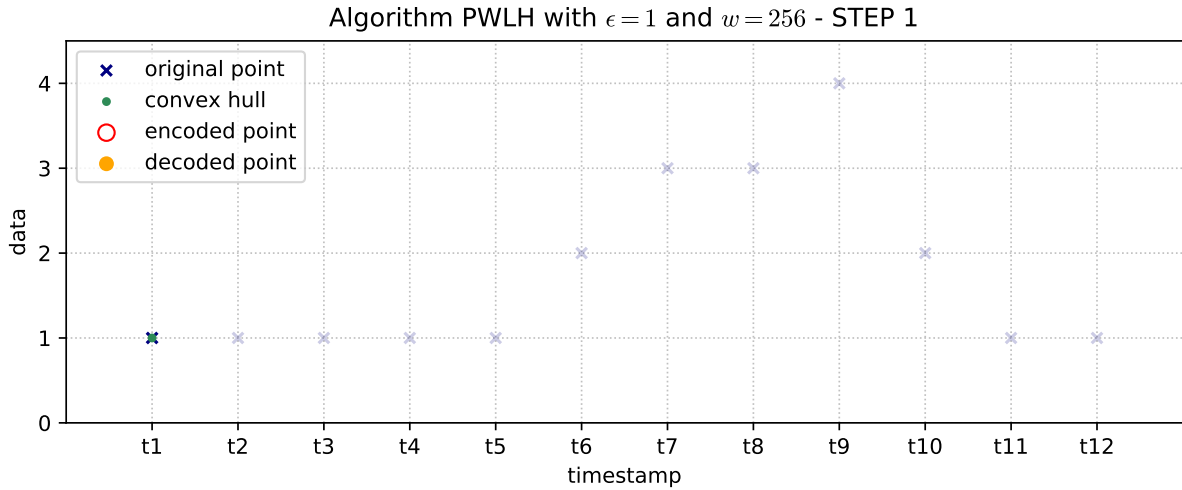


FIGURE 2.21

After adding the second data value to the window, the convex hull is updated. This step is shown in Figure 2.22. The convex hull consists of a single edge, and the maximum distance from either point to the edge is zero, i.e.  $width = 0 \leq 2 * \epsilon = 2$ , so the condition in line 14 is false and the window is not coded.

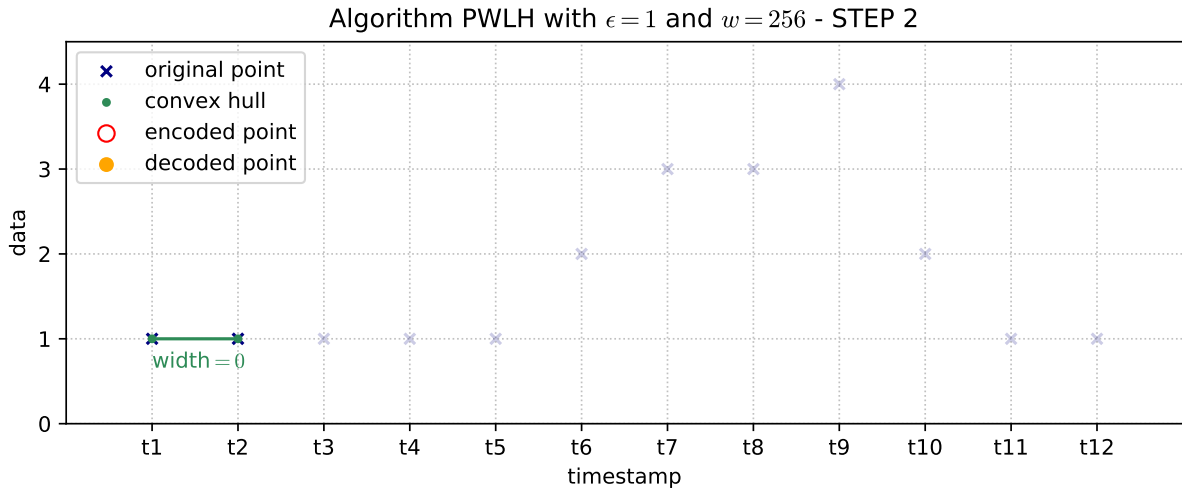


FIGURE 2.22

The next three data values are also equal to 1. They are added to the window, and the convex hull is updated, but *width* doesn't change in any case, so the window is not coded. This step is shown in Figure 2.23.

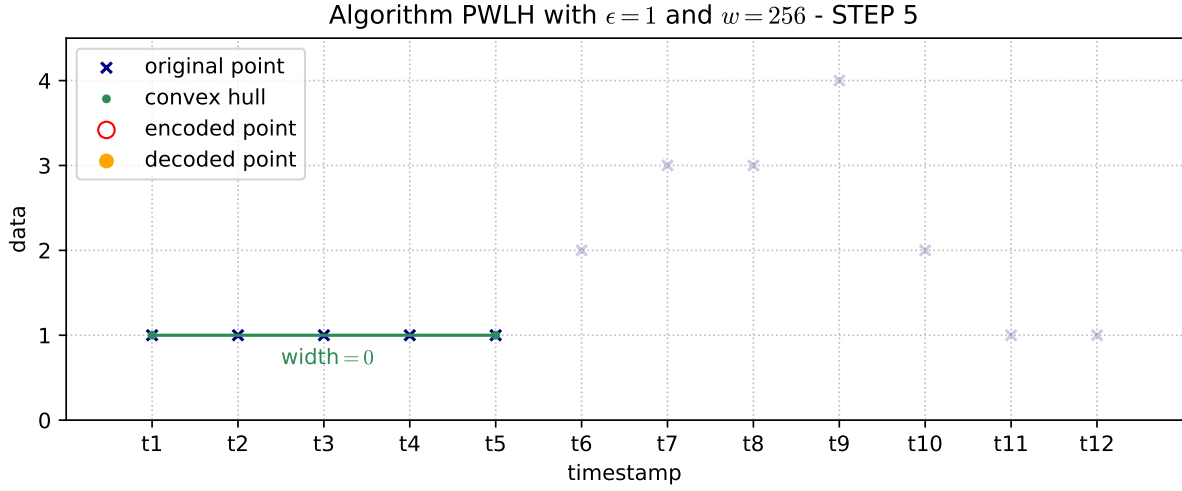


FIGURE 2.23

Next, the data value 2 is added to the window. The updated convex hull, which now consists of three edges, is shown in Figure 2.24. In this case, the maximum distance between the upper edge and any of the points in the convex hull is approximately 0.8. Therefore,  $width \approx 0.8 \leq 2$ , the condition in line 14 remains false, and the window is not coded.

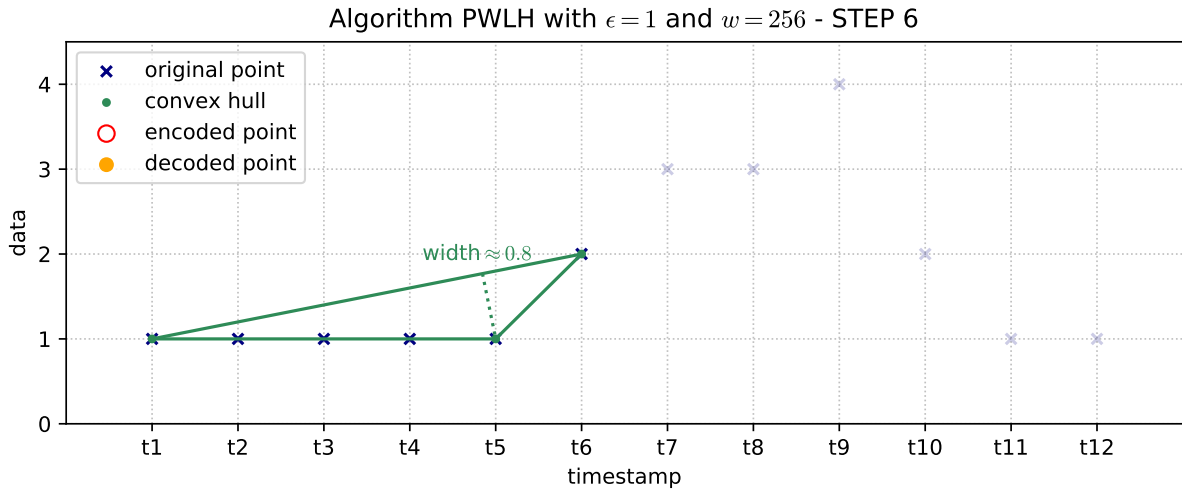


FIGURE 2.24



The following three iterations are quite similar to the previous one. In every case, the convex hull is updated, and even though the maximum distance between the upper edge and any of the points in the convex hull increases, it is never larger than 2, so the window is not coded. These steps are shown in figures 2.25, 2.26 and 2.27.

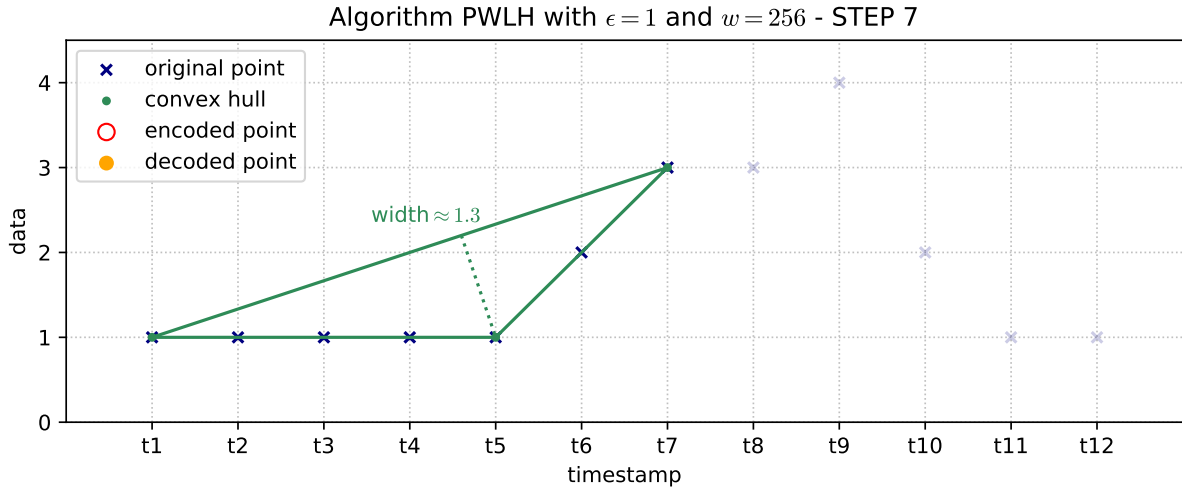


FIGURE 2.25

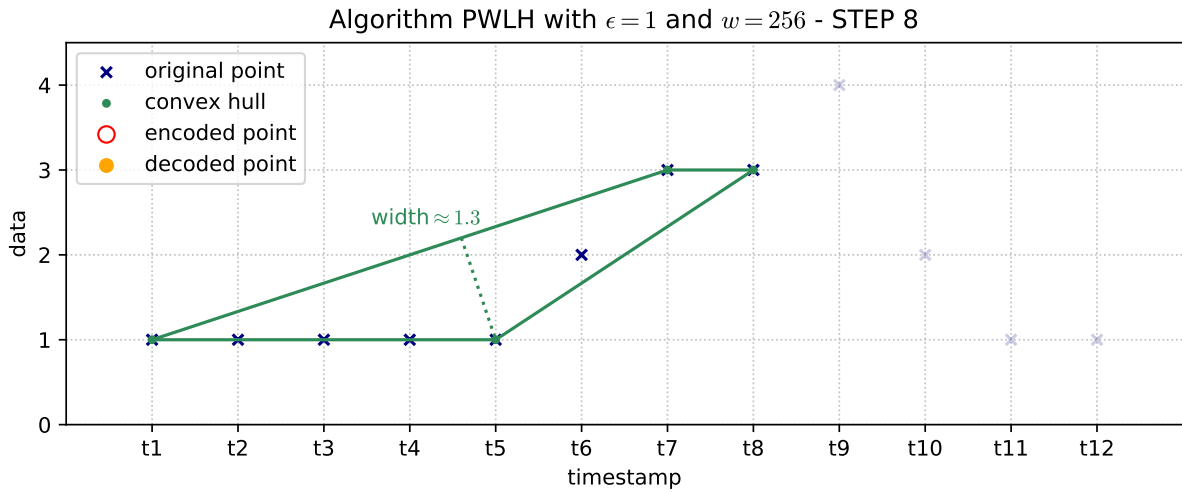


FIGURE 2.26

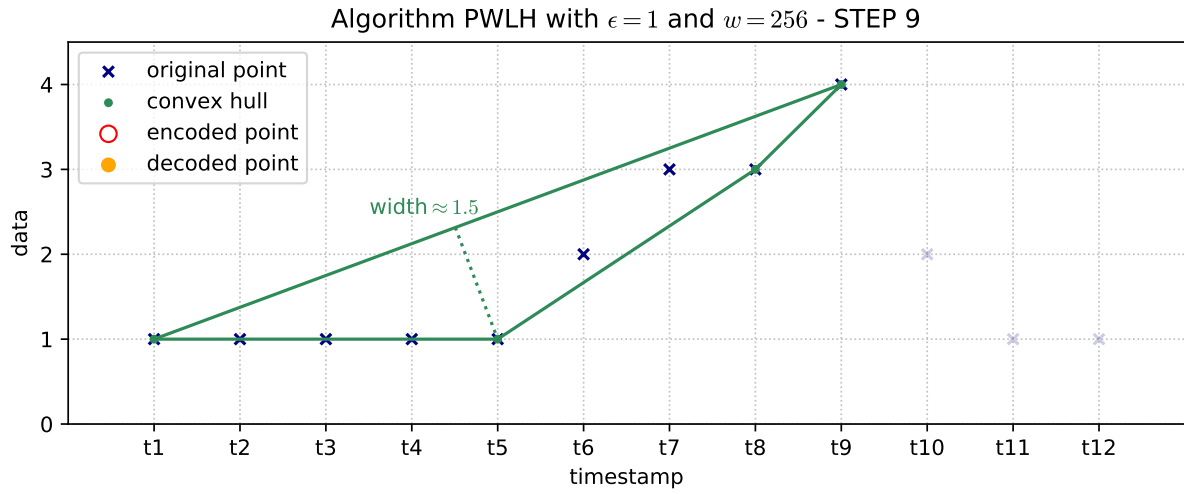


FIGURE 2.27

Eventually, in the 10th iteration, after adding the data value 2 to the window, the resulting convex hull, which is shown in Figure 2.28, violates the threshold condition for the first time. Observe that for every edge in the convex hull there exists a point in the hull such that its distance to the edge is larger than 2.

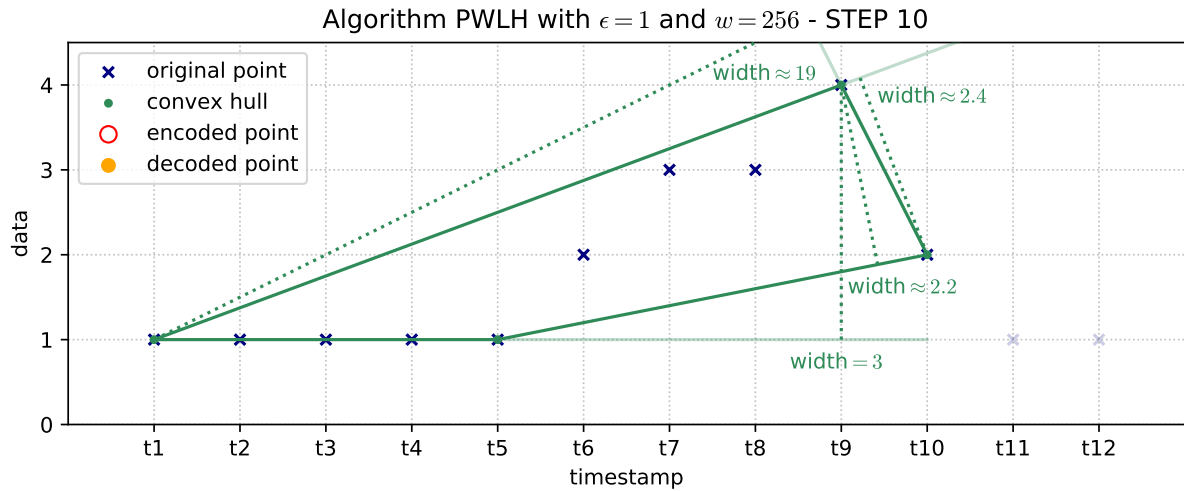


FIGURE 2.28

Since the condition in line 14 becomes true, the last entry is removed from the window and the convex hull is updated (lines 15-16), the window is encoded (lines 21-24), and the data value which violated the threshold condition is added to a new window and its convex hull (lines 25-27). Figure 2.29 shows the encoded values and the new convex hull, which consists of a single point. Notice that the line segment in the graph, whose endpoints were encoded, is the one that minimizes the maximum distance to the nine points in the encoded window.

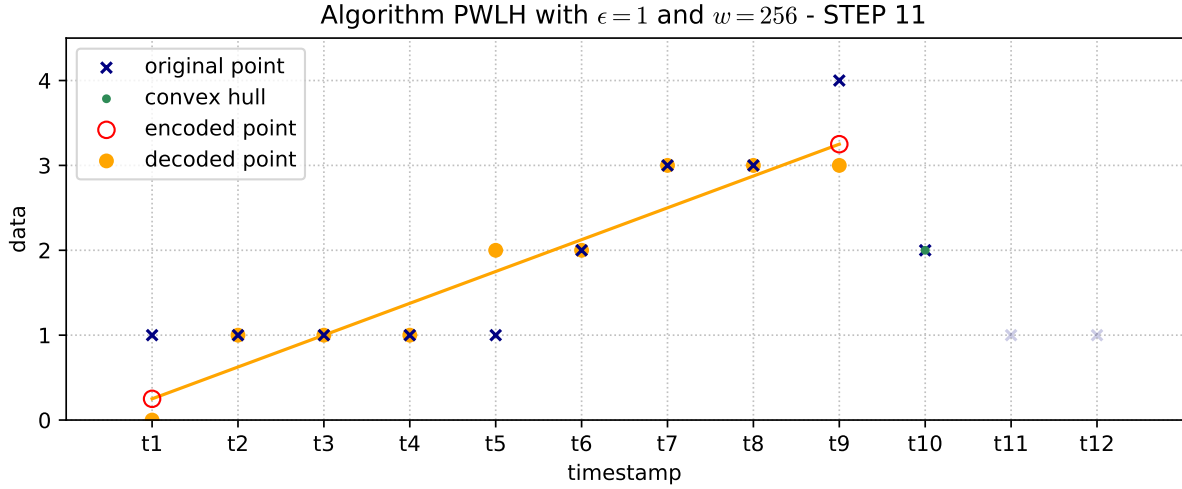


FIGURE 2.29

In the last two iterations, which correspond to the last two data values, the threshold condition is not violated. Therefore, after executing the last iteration, *win* is not empty. In this case, the algorithm must still encode its values, which is done by executing the same code as in lines 21-27. This was left out of the pseudocode for clarity. Figure 2.30 shows the line segment that minimizes the maximum distance to the three points in the second window.

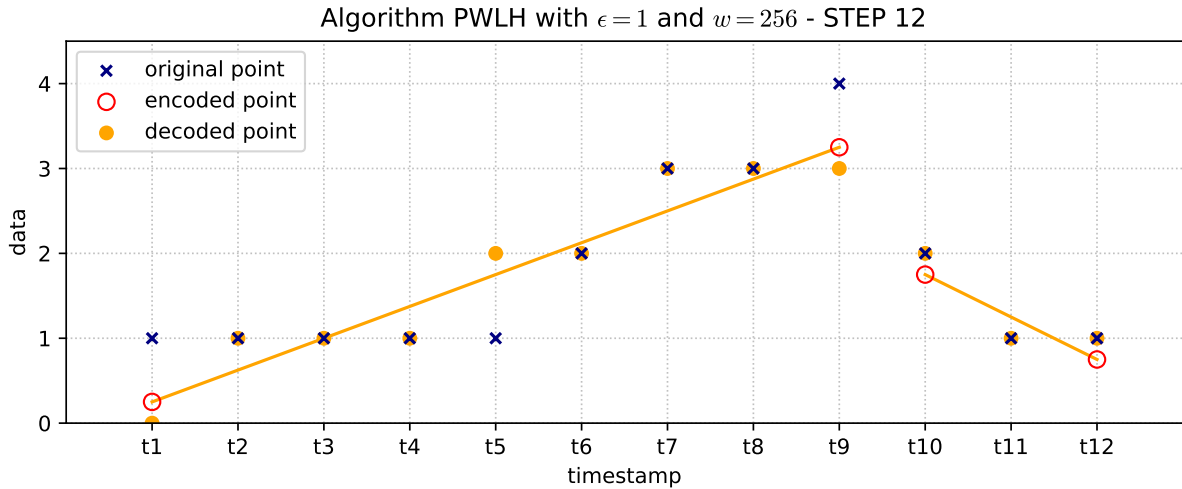


FIGURE 2.30

TODO: mencionar PWLHInt y explicar las tres diferencias con PWLH:

- (1) distinto  $\epsilon$
- (2) chequeo de que el endpoint esté dentro del rango de valores
- (2) `code_base_2` en vez de `code_float` (lines 23-24).

## 2.7 Algorithm SF

Es similar a PWLH (ver paper An Evaluation...).

Details are omitted...

Describe example in Figure 2.31.

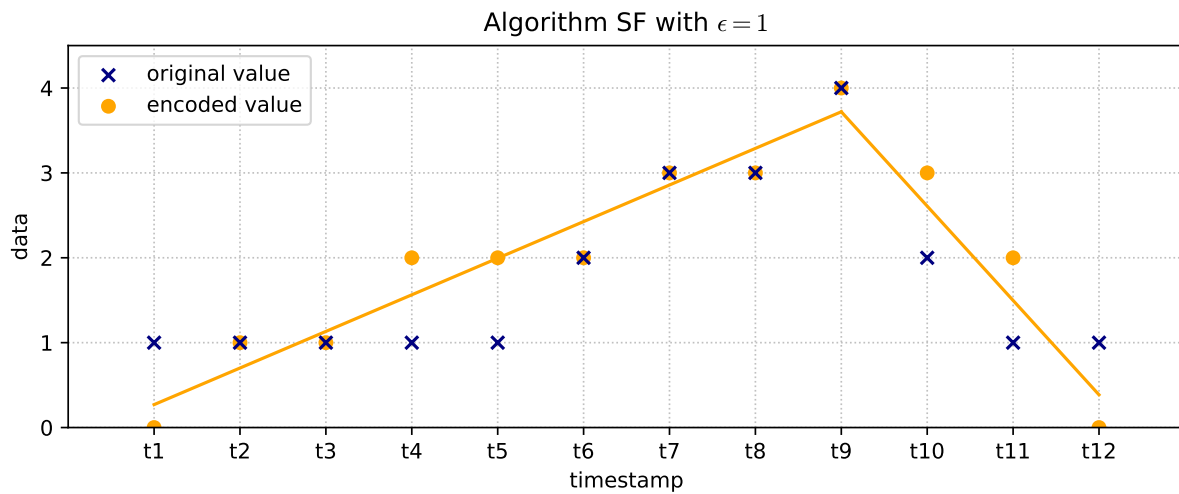


FIGURE 2.31: Example SF

## 2.8 Algorithm CA

Algorithm CA [12], also known as Critical Aperture, supports lossless and lossy compression, with both variants ( $M$  and  $NM$ ), and it has a window size parameter ( $w$ ) that establishes the maximum block size in which the data are processed and encoded. It is a Linear model algorithm, so it encodes signals using linear functions.

Figure 2.32 shows the pseudocode for the `code_column_M` subroutine for algorithm CA. TODO: continue... IDEA: Since the condition involves geometrical concepts it is further developed in the example...

```

input : column: column of the csv data file to be coded
        out: binary file encoded with algorithm CA
         $\epsilon$ : maximum error threshold
         $w$ : maximum window size
1  win = new_window()
2  foreach entry in column.entries do
3      code_window = false
4      code_value = false
5      if entry == column.entries[0] then
6          | code_value = true
7      else if win.size == 0 then
8          | win.add_incoming_point(entry)
9          | continue
10     else if win.size ==  $w$  or not win.CA_condition_holds?(entry,  $\epsilon$ ) then
11         | code_window = true
12         | code_value = true
13     end
14     if code_window then
15         | out.code_base_2(win.size - 1,  $\log_2 w$ )
16         | out.code_base_2(win.code_value, column.total_bits)
17     end
18     if code_value then
19         | out.code_base_2(0,  $\log_2 w$ )
20         | out.code_base_2(entry.value, column.total_bits)
21         | win.set_archived_point(entry)
22     end
23 end

```

FIGURE 2.32: CA.code\_column\_M pseudocode.

TODO:

- describir el ejemplo de las próximas 8 imágenes.
- mencionar que se tienen en cuenta los delta pero se omite en el pseudocódigo.
- mencionar que se tiene en cuenta el caso  $\text{delta} = 0 \Rightarrow$  por esto es que en el STEP 11 el valor de  $t_7$  se codifica aparte (si  $\text{delta} = 0$ , entonces  $t_6 = t_7$ , pero no necesariamente  $f(t_6) = f(t_7)$ )

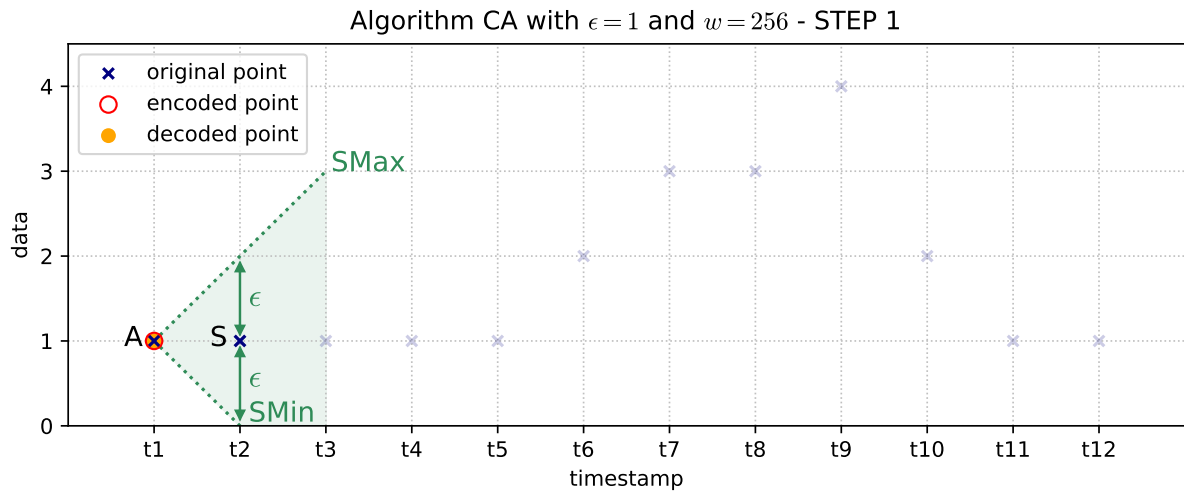


FIGURE 2.33

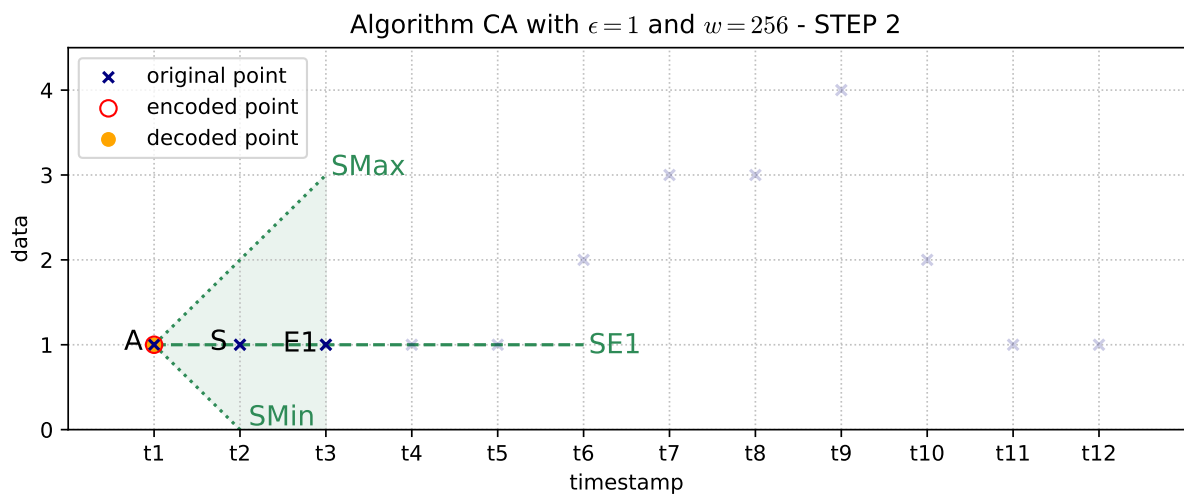


FIGURE 2.34

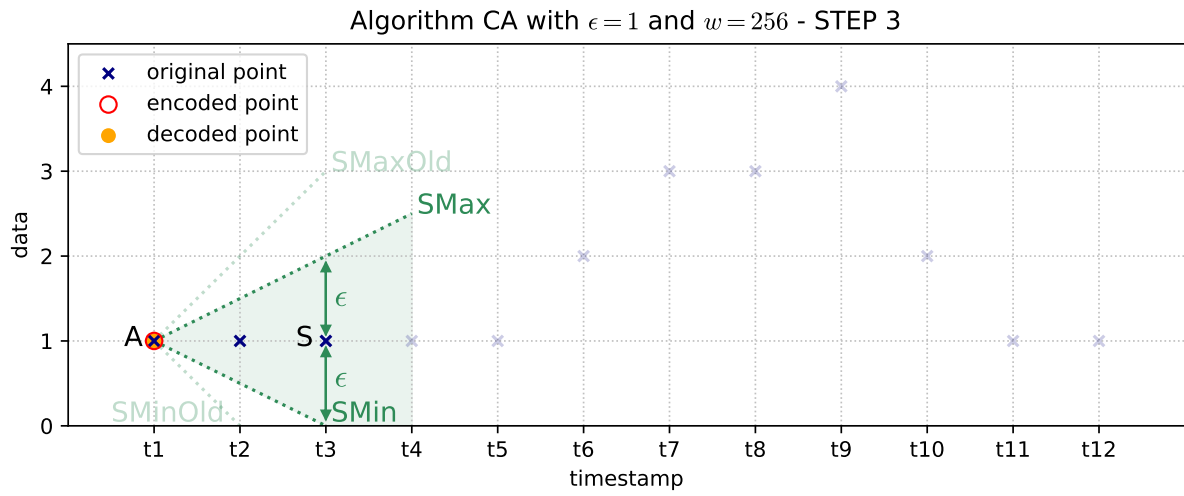


FIGURE 2.35

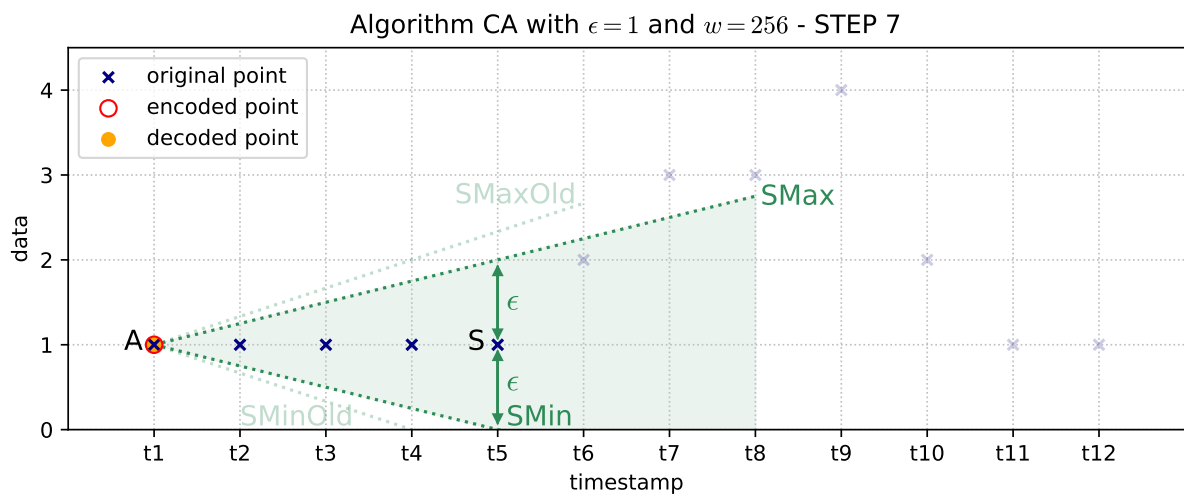


FIGURE 2.36



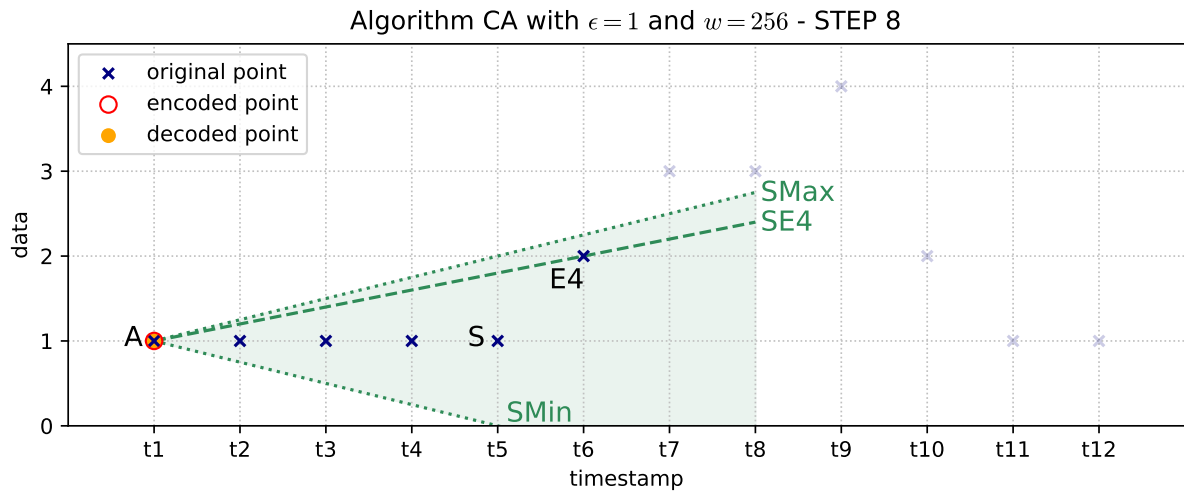


FIGURE 2.37

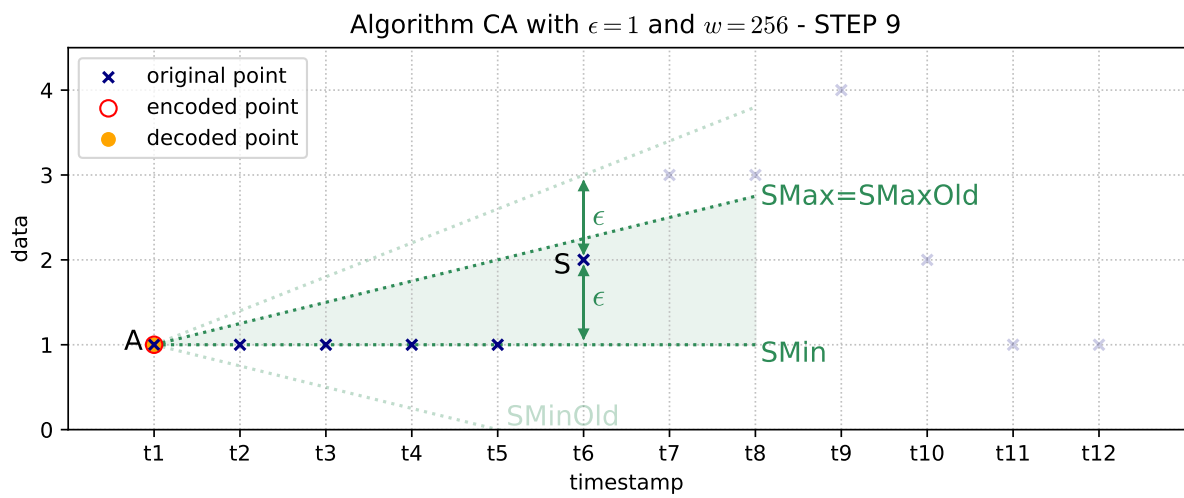


FIGURE 2.38

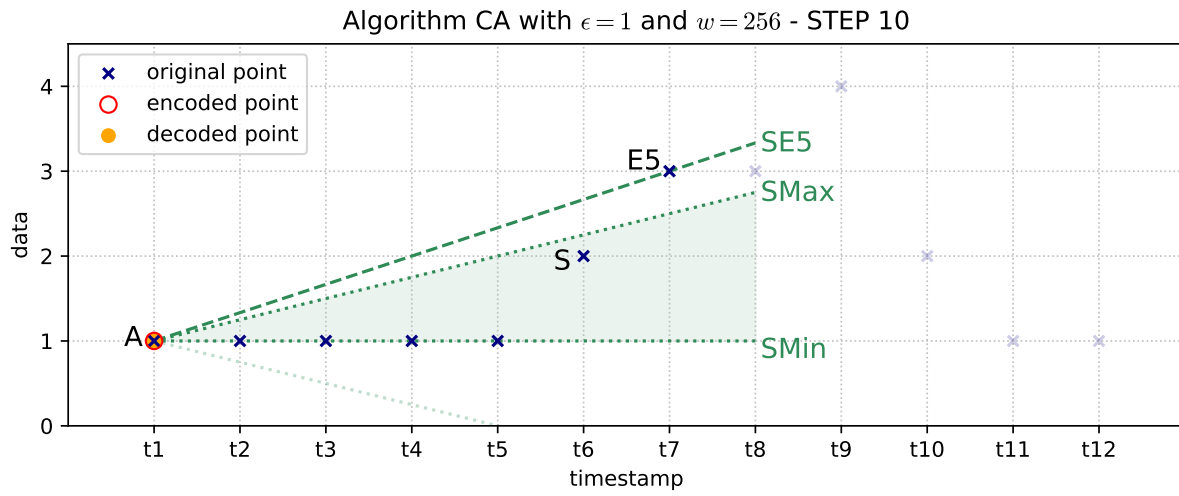


FIGURE 2.39

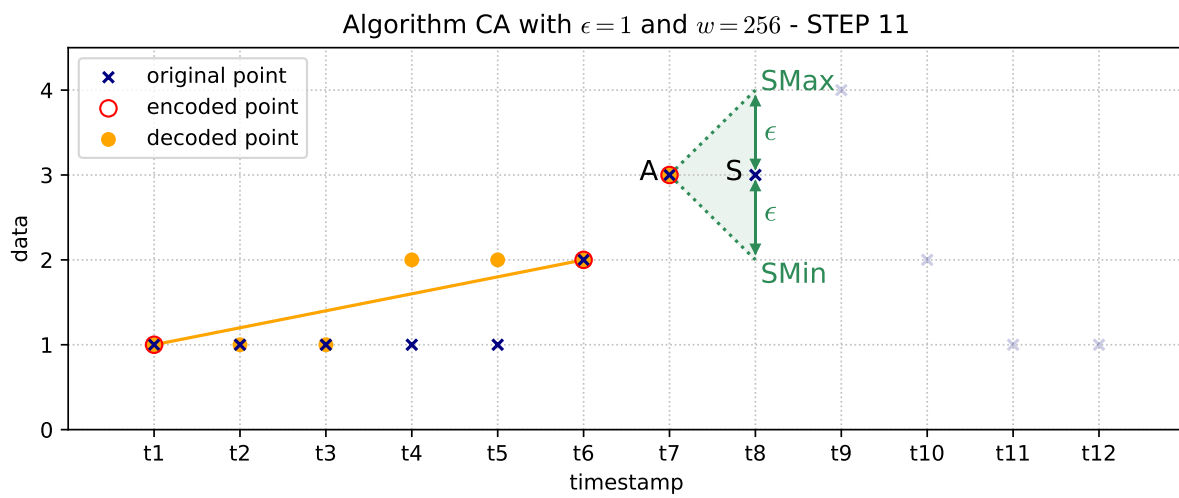


FIGURE 2.40

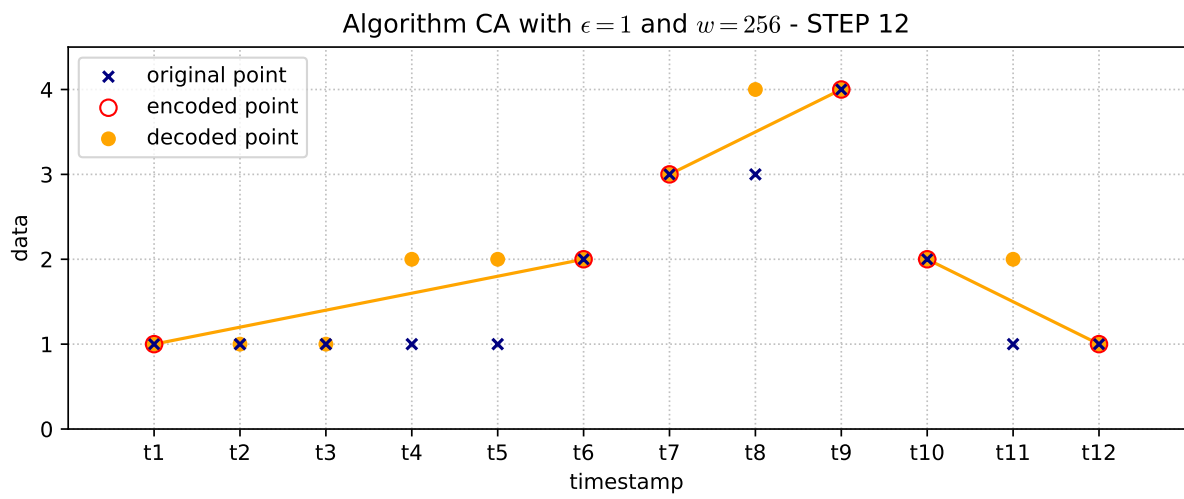


FIGURE 2.41

## 2.9 Algorithm FR

Algorithm FR [13], also known as Fractal Resampling, is an algorithm developed by the European Space Association (ESA). Since it was designed to be run on spacecraft and planetary probes, its computational complexity is fairly low. It supports lossless and lossy compression, with both variants ( $M$  and  $NM$ ), and it has a window size parameter ( $w$ ) that establishes the maximum block size in which the data are processed and encoded. It is a Linear model algorithm, so it encodes signals using linear functions. In particular, the data points in each window are modeled by one or more line segments that are selected using a simple recursive technique called mid-point displacement.

In Figure 2.42 we show the pseudocode for the `code_column_M` subroutine. Since this subroutine corresponds to the masking variant, we can assume that all of the processed column entries are integer values. We observe that these column entries are added to the window (line 3) until its size is equal to  $w$  (i.e. the condition in line 4 becomes false). When that happens, an empty array called `points_indexes` is created (line 7), and then passed as an argument to the `push_points_indexes` method (line 8), which fills it with integers in the range  $[0, w - 1]$  that correspond to the indexes of certain entries of `win`, the current window. The pseudocode with the details of that recursive method is presented later, but the key idea is that the points in `win` can be modeled by one or more line segments, and the `win` indexes of said line segment's endpoints are included, in order, in the `points_indexes` array. This array always includes the first and the last index of `win`, which implies that the first line segment always starts in the first point of the window, and the last line segment always ends in the last point of the window. In lines 9-12, every one of the endpoints is encoded, each using a total of  $\log_2 w + \text{column.total\_bits}$  bits. Observe that  $\log_2 w$  bits are required for encoding the window index of an endpoint, since said index is guaranteed to be a value between 0 and  $w - 1$ . Finally, a new empty window is created (line 13), and the next column entry is processed.

```

input : column: column of the csv data file to be coded
         out: binary file encoded with algorithm FR
          $\epsilon$ : maximum error threshold
          $w$ : maximum window size

1 win = new_window()
2 foreach entry in column.entries do
3   win.push(entry)
4   if win.size <  $w$  then
5     continue
6   end
7   points_indexes = new_array()
8   push_points_indexes(win,  $\epsilon$ , points_indexes, 0, win.size - 1)
9   foreach index in points_indexes do
10    out.code_base_2(index,  $\log_2 w$ )
11    out.code_base_2(win.entries[index].value, column.total_bits)
12  end
13  win = new_window()
14 end

```

FIGURE 2.42: FR.code\_column\_M pseudocode.

It should be pointed out that, after executing the last iteration, `win` might not be empty. In that case, the algorithm must still encode its values, which is done by executing the same code as in lines 7-12. This was left out of the pseudocode for clarity.

In Figure 2.43 we present the pseudocode for the `decode_column_M` subroutine. This subroutine keeps running until every entry in the column has been decoded. In each iteration, first a window with exactly  $w$  entries with empty values is created (line 2). Next there's a loop in which both the window index and the value corresponding to each segment endpoint are decoded (lines 4-5) and added to the window (line 6). The algorithm breaks out of the loop once the last endpoint is decoded (lines 7-9). Finally, the algorithm iterates through the window entries, writing each decoded value into the csv data file. If the value was directly read from the binary file (i.e. the condition in line 12 is true), then it corresponds to a segment endpoint, so no additional calculations are needed. Otherwise, the value can be obtained by using the index to replace the  $x$  variable in the equation of the line segment that corresponds to said index (line 15).

```

input : in: binary file coded with algorithm FR
         out: decoded csv data file
         w: maximum window size
1 while not in.column_decoded? do
2   win = new_window(w)
3   while true do
4     index = in.decode_base_2(log2 w)
5     value = in.decode_base_2(column.total_bits)
6     win.entries[index].value = value
7     if index == w - 1 then
8       break
9     end
10  end
11  foreach entry in win.entries do
12    if entry.value then
13      value = entry.value
14    else
15      value = win.segment_equation(entry.index)
16    end
17    out.write_string(value)
18  end
19 end

```

FIGURE 2.43: FR.decode\_column\_M pseudocode.

In Figure 2.44 we present the pseudocode for the `push_points_indexes` recursive method. As we recall, it is invoked from the `code_column_M` subroutine (line 8), and it fills the `points_indexes` array with the window indexes of the endpoints of one or more line segments. Since encoding each of these endpoints requires additional bits, the algorithm must try to minimize the amount of entries in the array, while avoiding the error threshold constraint being violated.

Next we present an example that illustrates the mechanics of algorithm FR and its main method, `push_points_indexes`. Again, we describe the process of encoding the data values which are shown in Figure 2.8, and the comments regarding the timestamp values and the graph legends that were made when introducing said figure, also apply in this case. In this example we use algorithm FR with an error threshold parameter ( $\epsilon$ ) equal to 1, and a maximum window size ( $w$ ) equal to 256.

The condition in line 4 of the `code_column_M` subroutine is true for every iteration, so every data value, i.e.  $[1, 1, 1, 1, 1, 2, 3, 3, 4, 2, 1, 1]$ , is added to the window. After executing the last iteration, *win* is not empty, so as we pointed out, the algorithm must encode its values by executing the same code as in lines 7-12. The `push_points_indexes` method is first called with parameters *win*,  $\epsilon = 1$ , `points_indexes` = [], *first\_index* = 0 and *first\_index* = 11. After executing lines 1-6, we have `points_indexes` = [0, 11]. This step is shown in Figure 2.45.

```

input : win: window
         $\epsilon$ : maximum error threshold
        points_indexes: array with the win index of the displaced points
        first_index: index of the first win entry to be processed
        last_index: index of the last win entry to be processed
1 if not points_indexes.includes?(first_index) then
2   | points_indexes.push_and_sort(first_index)
3 end
4 if not points_indexes.includes?(last_index) then
5   | points_indexes.push_and_sort(last_index)
6 end
7 if first_index + 1 < last_index and
   not win.FR_condition_holds?(first_index, last_index,  $\epsilon$ ) then
8   | half = (first_index + last_index)/2
9   | push_points_indexes(win,  $\epsilon$ , points_indexes, first_index, half)
10  | push_points_indexes(win,  $\epsilon$ , points_indexes, half, last_index)
11 end

```

FIGURE 2.44: push\_points\_indexes pseudocode.

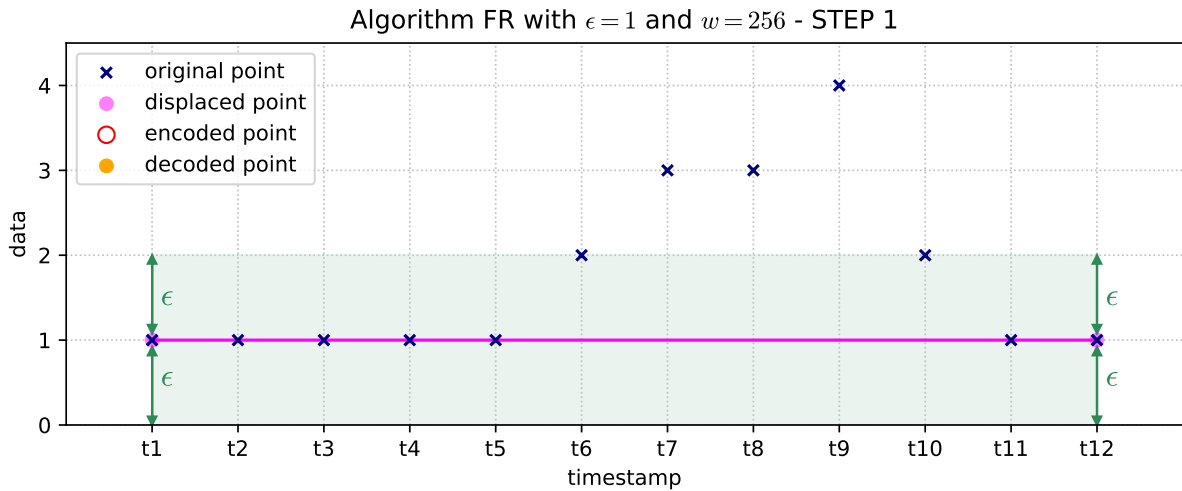


FIGURE 2.45

The condition in line 7 is true, since  $0 + 1 < 11$ . The condition in line 8 is also true, since as Figure 2.45 shows, the points the corresponding timestamps  $t_7$ ,  $t_8$ , and  $t_9$ , violate the error threshold constraint. Since both conditions are true, the `push_points_indexes` method is called twice (lines 9-10).

The first time it is called with parameters `points_indexes = [0, 11]`, `first_index = 0` and `last_index = 5` and it adds a single index, making `points_indexes = [0, 5, 11]`. This step is shown in Figure 2.46. Since the threshold constraint is satisfied between  $t_1$  and  $t_6$ , the condition in line 8 is false, so no further calls to the `push_points_indexes` are made.

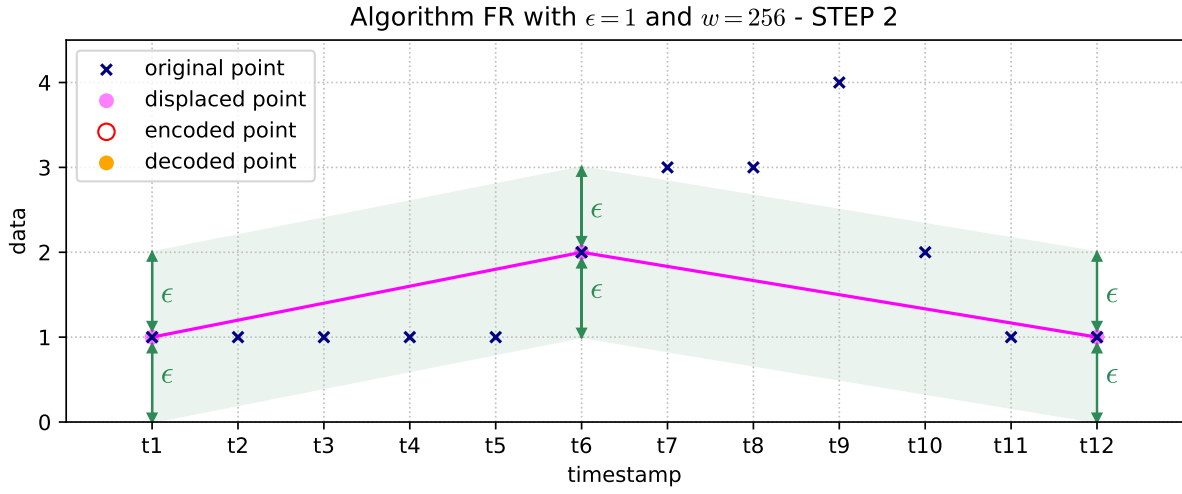


FIGURE 2.46

The second time it is called with parameters  $points\_indexes = [0, 5, 11]$ ,  $first\_index = 5$  and  $last\_index = 11$ . In this case, as Figure 2.46 shows, the points corresponding to timestamps  $t_7$ ,  $t_8$  and  $t_9$  still violate the error threshold constraint, which means that, once again, the `push_points_indexes` method must be called twice.

After the first invocation of the `push_points_indexes` method is complete (line 8 of the code `_column_M` subroutine), we have  $points\_indexes = [0, 5, 8, 11]$ . This information is shown in Figure 2.47. Observe that all of the points satisfy the error threshold constraint.

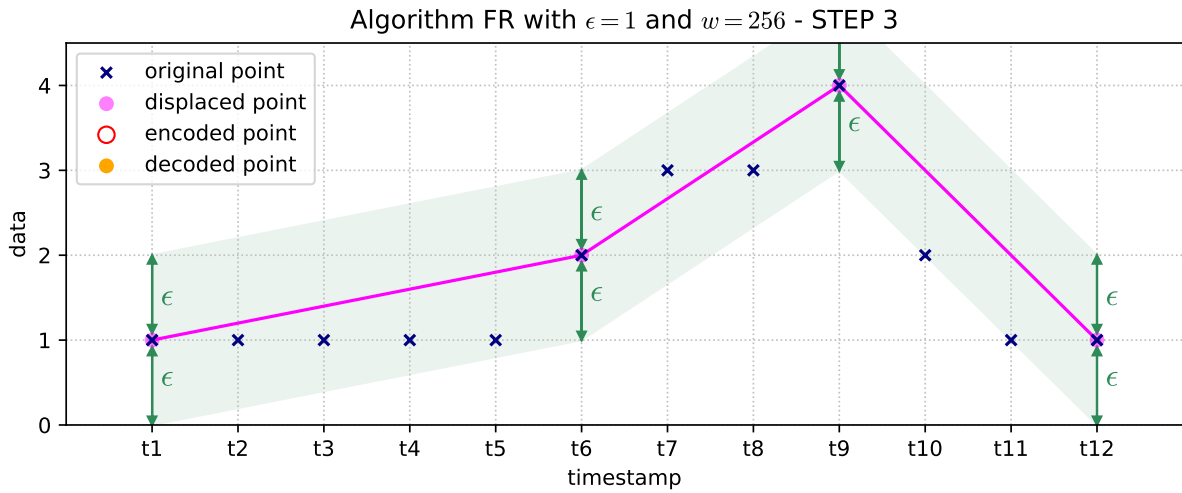


FIGURE 2.47

Finally, in lines 9-12 of the `_column_M` subroutine, the four endpoints associated to the indexes in  $points\_indexes$  are encoded. The graph in Figure 2.48 shows the original and the encoded data values, and the three segments associated to the four endpoints. As we recall from the pseudocode in Figure 2.43, the `decode_column_M` subroutine reads the values corresponding to the segment endpoints directly from the coded binary file, so in these cases the original and the encoded values always match.

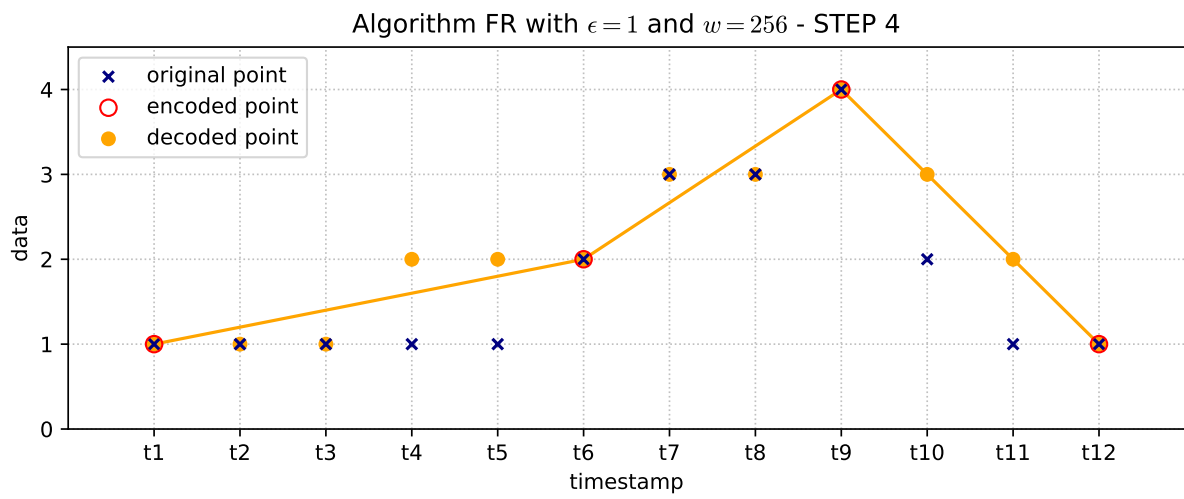


FIGURE 2.48

Cuando se hacen las proyecciones hay que tener en cuenta la time column...



## 2.10 Algorithm GAMPS

Ver los siguientes documentos:

- [08] [AVANCES / DUDAS](#)
- [09] [AVANCES / DUDAS](#)
- [10] [AVANCES / DUDAS](#)
- [11] [AVANCES / DUDAS](#)
- [12] [AVANCES / DUDAS](#)

## Chapter 3

# Experimental Results

In this chapter we present our experimental results. The main goal of our experiments is to analyze the performance of each of the coding algorithms presented in Chapter 2, by encoding the various datasets introduced in Chapter 1.

In Section 3.1 we describe our experimental setting and define the evaluated combinations of algorithms, their variants and parameter values, and the figures of merit used for comparison.

In Section 3.2 we compare the compression performance of the masking and non-masking variants for each coding algorithm. The results show that on datasets with few or no gaps the performance of both variants is roughly the same, while on datasets with many gaps the masking variant always performs better, in some cases with a significative difference. These results suggest that the masking variant is more robust and performs better in general.

In Section 3.3 we analyze the extent to which the window size parameter **impacts/affects** the compression performance of the coding algorithms. We compress each dataset file, and compare the results obtained when using the optimal window size (i.e. the one that achieves the best compression) for each file, with the results obtained when using the optimal window size for the whole dataset. The results indicate that the **impact/effect** of using the optimal window size for the whole dataset instead of the optimal window size for each file is rather small.

In Section 3.4 we compare the performance of the different coding algorithms among each other and with the general purpose compression algorithm gzip. Among the tested coding algorithms, for larger error thresholds APCA is the best algorithm for compressing every data type in our experimental data set, while for lower thresholds the recommended algorithms are PCA, APCA and FR, depending on the data type. If we also consider algorithm gzip, **there isn't an algorithm that is better for compressing every data type for any threshold value / for no threshold value there exists an algorithm that is better for compressing every data type**. Depending on the data type, the recommended algorithms are APCA and gzip for larger error thresholds, and PCA, APCA, FR and gzip for lower thresholds.

### 3.1 Experimental Setting

We denote by  $A$  the set of all the coding algorithms presented in Chapter 2. For an algorithm  $a \in A$ , we denote by  $a_v$  its variant  $v$ , where  $v$  can be  $M$  (masking) or  $NM$  (non-masking). There exist some  $a \in A$  for which either  $a_M$  or  $a_{NM}$  is invalid (recall this information from Table 2.1). We denote by  $V$  the set of variants composed of every valid variant  $a_v$  for every algorithm  $a \in A$ . Also, we denote by  $A_M$  the subset of algorithms from  $A$  composed of every algorithm for which both variants,  $a_M$  and  $a_{NM}$ , are valid.

We evaluate the compression performance of every algorithm  $a \in A$  on the datasets described in Chapter 1. For each algorithm we test every valid variant  $a_v$ . We also test several combinations of algorithm parameters. Specifically, for the algorithms that admit a window size parameter  $w$  (every algorithm except *Base* and *SF*), we test all the values of  $w$  in the set  $W = \{4, 8, 16, 32, 64, 128, 256\}$ . For the encoders that admit a lossy compression mode with a threshold parameter  $e$  (every encoder except *Base*), we test all the values of  $e$  in the set  $E = \{1, 3, 5, 10, 15, 20, 30\}$ , where each threshold is expressed as a percentage fraction of the standard deviation of the data being encoded. For example, for certain data with a standard deviation of 20, taking  $e = 10$  implies that the lossy compression allows for a maximal per-sample distortion of 2 sampling units.

**Definition 3.1.1.** We refer to a specific combination of a coding algorithm variant and its parameter values as a *coding algorithm instance (CAI)*. We define  $CI$  as the set of all the CAIs obtained by combining each of the variants  $a_v \in V$  with the parameter values (from  $W$  and  $E$ ) that are suitable for algorithm  $a$ . We denote by  $c_{\langle a_v, w, e \rangle}$  the CAI obtained by setting a window size parameter equal to  $w$  and a threshold parameter equal to  $e$  on algorithm variant  $a_v$ .

We assess the compression performance of a CAI mainly through the compression ratio, which we define next. For this definition, we regard *Base* as a trivial CAI that serves as a base ground for compression performance comparison (recall the definition of algorithm *Base* from Section 2.3).

**Definition 3.1.2.** Let  $f$  be a file and  $z$  a data type of a certain dataset. We define  $f_z$  as the subset of data of type  $z$  from file  $f$ . For example, for the dataset Hail, the data type  $z$  may be Latitude, Longitude, or Size.

**Definition 3.1.3.** Let  $f$  be a file and  $z$  a data type of a certain dataset. Let  $c \in CI$  be a CAI. We define  $|c(z, f)|$  as the size in bits of the resulting bit stream obtained by coding  $f_z$  with  $c$ .

**Definition 3.1.4.** The *compression ratio (CR)* of a CAI  $c \in CI$  for the data type  $z$  of a certain file  $f$  is the fraction of  $|c(z, f)|$  with respect to  $|\text{Base}(z, f)|$ , i.e.,

$$CR(c, z, f) = \frac{|c(z, f)|}{|\text{Base}(z, f)|}. \quad (3.1)$$

Notice that smaller values of CR correspond to better performance. Our main goals are to analyze which CAIs yield the smallest values in (3.1) for the different data types, and to study how the CR depends on the different algorithms, their variants and the parameter values.

To compare the compression performance between a pair of CAIs we calculate the relative difference, which we define next.

**Definition 3.1.5.** The *relative difference* ( $RD$ ) between a pair of CAIs  $c_1, c_2 \in CI$  for the data type  $z$  of a certain file  $f$  is given by

$$RD(c_1, c_2, z, f) = 100 \times \frac{|c_2(z, f)| - |c_1(z, f)|}{|c_2(z, f)|}. \quad (3.2)$$

Notice that  $c_1$  has a better performance than  $c_2$  if (3.2) is positive.

In some of our experiments we consider the performance of algorithms on complete datasets, rather than individual files. With this in mind, we extend the definitions 3.1.3–3.1.5 to datasets, as follows.

**Definition 3.1.6.** Let  $z$  be a data type of a certain dataset  $d$ . We define  $F(d, z)$  as the set of files  $f$  from dataset  $d$  for which  $f_z$  is not empty.

**Definition 3.1.7.** Let  $z$  be a data type of a certain dataset  $d$ . Let  $c \in CI$  be a CAI. We define  $|c(z, d)|$  as

$$|c(z, d)| = \sum_{f \in F(d, z)} |c(z, f)|. \quad (3.3)$$

**Definition 3.1.8.** The *compression ratio* ( $CR$ ) of a CAI  $c \in CI$  for the data type  $z$  of a certain dataset  $d$  is given by

$$CR(c, z, d) = \frac{|c(z, d)|}{|\text{Base}(z, d)|}. \quad (3.4)$$

**Definition 3.1.9.** The *relative difference* ( $RD$ ) between a pair of CAIs  $c_1, c_2 \in CI$  for the data type  $z$  of a certain dataset  $d$  is given by

$$RD(c_1, c_2, z, d) = 100 \times \frac{|c_2(z, d)| - |c_1(z, d)|}{|c_2(z, d)|}. \quad (3.5)$$

## 3.2 Comparison of Masking and Non-Masking Variants

In this section, we compare the compression performance of the masking and non-masking variants of every coding algorithm in  $A_M$ . Specifically, we compare:

- $\text{PCA}_M$  against  $\text{PCA}_{NM}$
- $\text{APCA}_M$  against  $\text{APCA}_{NM}$
- $\text{CA}_M$  against  $\text{CA}_{NM}$
- $\text{PWLH}_M$  against  $\text{PWLH}_{NM}$
- $\text{PWLHInt}_M$  against  $\text{PWLHInt}_{NM}$
- $\text{GAMPSLimit}_M$  against  $\text{GAMPSLimit}_{NM}$

For each algorithm  $a \in A_M$  and each threshold parameter, we compare the performance of  $a_M$  and  $a_{NM}$ . For the purpose of this comparison, we choose the most favorable window size for each variant  $a_v$ , in the sense of the following definition.

**Definition 3.2.1.** The *optimal window size (OWS)* of a coding algorithm variant  $a_v \in V$ , and a threshold parameter  $e \in E$ , for the data type  $z$  of a certain dataset  $d$ , is given by

$$\text{OWS}(a_v, e, z, d) = \arg \min_{w \in W} \left\{ CR(c_{<a_v, w, e>}, z, d) \right\}, \quad (3.6)$$

where we break ties in favor of the smallest window size.

For each data type  $z$  of each dataset  $d$ , and each coding algorithm  $a \in A_M$  and threshold parameter  $e \in E$ , we calculate the RD between  $c_{<a_M, w_M^*, e>}$  and  $c_{<a_{NM}, w_{NM}^*, e>}$ , as defined in (3.5), where  $w_M^* = \text{OWS}(a_M, e, z, d)$  and  $w_{NM}^* = \text{OWS}(a_{NM}, e, z, d)$ .

As an example, in figures 3.1 and 3.2 we show the CR and the RD, as a function of the threshold parameter, obtained for two data types of two different datasets. Figure 3.1 shows the results for the data type “SST” of the dataset SST, and Figure 3.2 shows the results for the data type “Longitude” of the dataset Tornado. In Figure 3.1 we observe a large RD favoring the masking variant for all tested algorithms. On the other hand, in Figure 3.2 we observe that the non-masking variant outperforms the masking variant for all algorithms. We notice, however, that the RD is very small in the latter case.

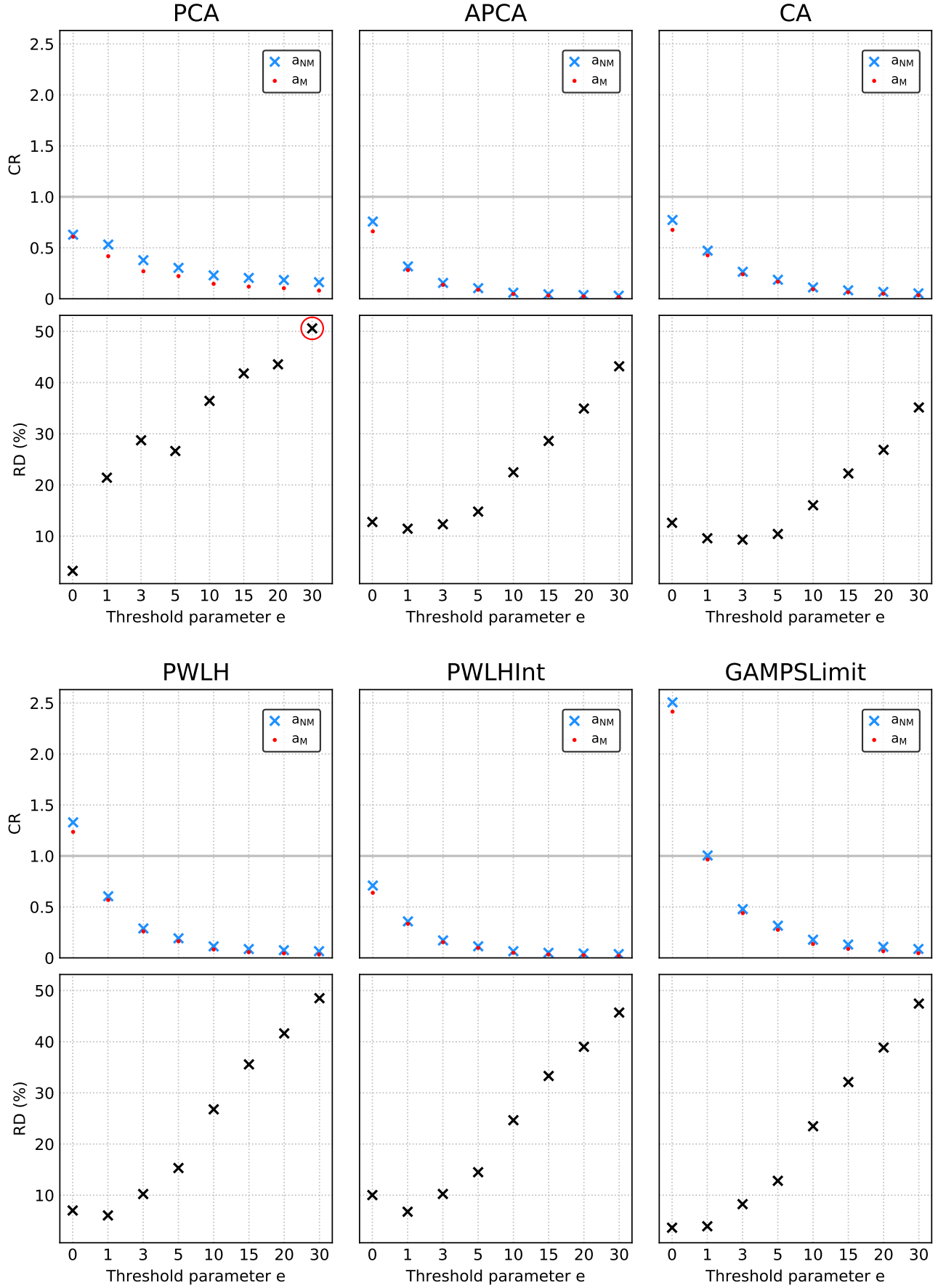


FIGURE 3.1: CR and RD plots for every pair of algorithm variants  $a_M, a_{NM} \in A_M$ , for the data type “SST” of the dataset SST. In the RD plot for algorithm PCA we highlight with a red circle the marker for the maximum value (50.60%) obtained for all the tested CAIs.

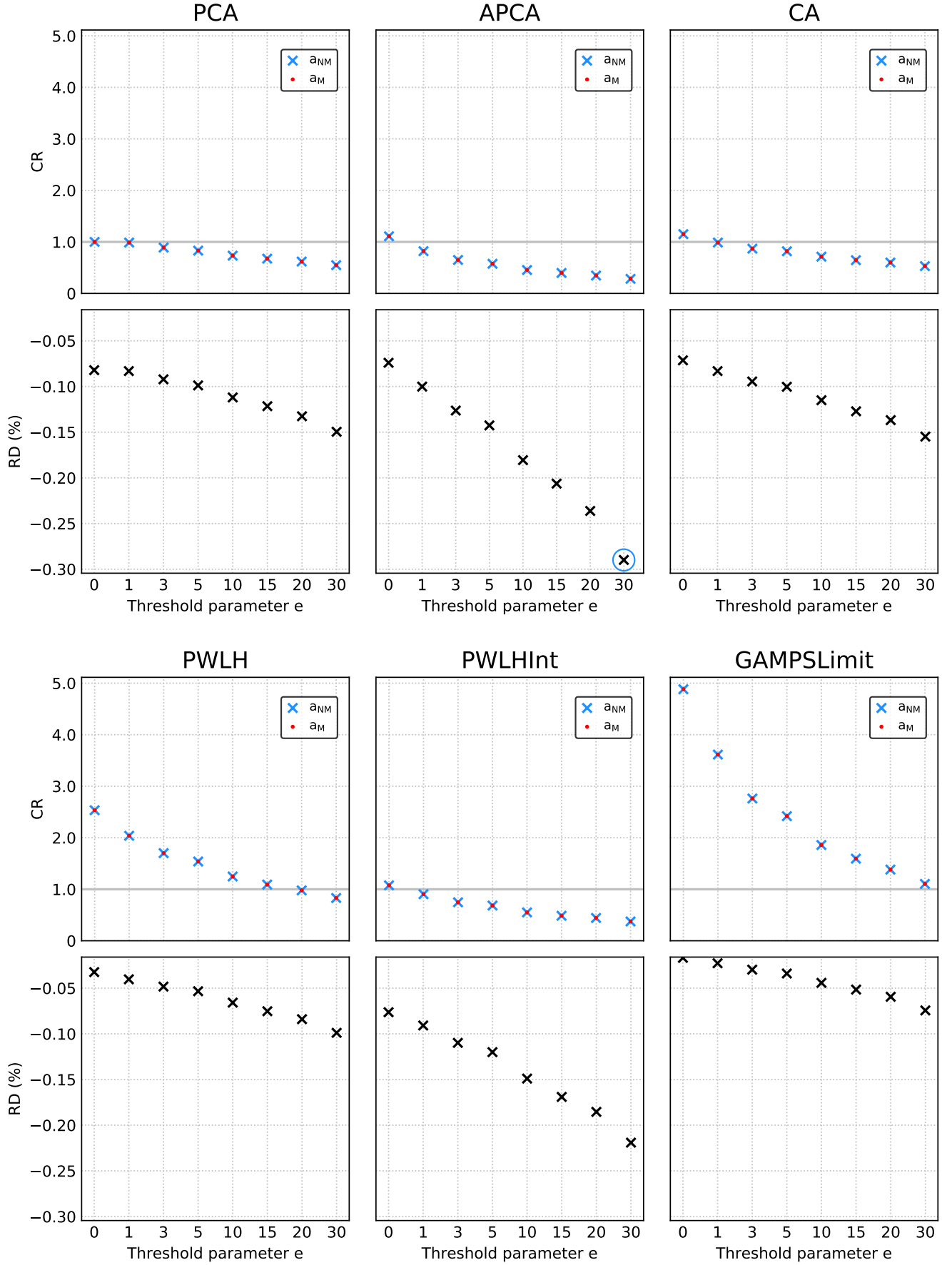


FIGURE 3.2: CR and RD plots for every pair of algorithm variants  $a_M, a_{NM} \in A_M$ , for the data type “Longitude” of the dataset Tornado. In the RD plot for algorithm APCA we highlight with a blue circle the marker for the minimum value (-0.29%) obtained for all the tested CAIs.

We analyze the experimental results to compare the performance of the masking and non-masking variants of each algorithm. For each data type, we iterate through each algorithm  $a \in A_M$ , and each threshold parameter  $e \in E$ , and we calculate the RD between the CAIs  $c_{\langle a_M, w_M^*, e \rangle}$  and  $c_{\langle a_{NM}, w_{NM}^*, e \rangle}$ , obtained by setting the OWS for the masking variant  $a_M$  and the non-masking variant  $a_{NM}$ , respectively. Since we consider 8 threshold parameters and there are 6 algorithms in  $A_M$ , for each data type we compare a total of 48 pairs of CAIs. Table 3.1 summarizes the results of these comparisons, aggregated by dataset. The number of pairs of CAIs evaluated for each dataset depends on the number of different data types it contains.

Dataset	Dataset Characteristic	Cases where $a_M$ outperforms $a_{NM}$ (%)	RD (%) Range
IRKIS	Many gaps	48/48 (100%)	(0; 36.88]
SST	Many gaps	48/48 (100%)	(0; 50.60]
ADCP	Many gaps	48/48 (100%)	(0; 17.35]
ElNino	Many gaps	336/336 (100%)	(0; 50.52]
Solar	Few gaps	73/144 (50.7%)	[-0.25; 1.77]
Hail	No gaps	0/144 (0%)	[-0.04; 0]
Tornado	No gaps	0/96 (0%)	[-0.29; 0]
Wind	No gaps	0/144 (0%)	[-0.12; 0]

TABLE 3.1: Range of values for the RD between the masking and non-masking variants of each algorithm (last column); we highlight the maximum (red) and minimum (blue) values taken by the RD. The results are aggregated by dataset. The second column indicates the characteristic of each dataset, in terms of the amount of gaps. The third column shows the number of cases in which the masking variant outperforms the non-masking variant of a coding algorithm, and its percentage among the total pairs of CAIs compared for a dataset.

Consider, for example, the results for the dataset Wind, in the last row. The second column shows that there are no gaps in any of the data types of the dataset (recall the dataset description from Table 1.1). Since the dataset has three data types, we compare a total of  $3 \times 48 = 144$  pairs of CAIs. The third column reveals that in none of these comparisons the masking variant  $a_M$  outperforms the non-masking variant  $a_{NM}$ , i.e. the RD is always negative. The last column shows the range for the values attained by the RD for those tested CAIs.

Observing the last column of Table 3.1, we notice that in every case in which the non-masking variant performs best, the RD is close to zero. The minimum value it takes is -0.29%, which is obtained for the data type “Longitude” of the dataset Tornado, with algorithm APCA, and error parameter  $e = 30$ . In Figure 3.2 we highlight the marker associated to this minimum with a blue circle. On the other hand, we also notice that for the datasets in which the masking variant performs best, the RD reaches high absolute values. The maximum (50.60%) is obtained for the data type “VWC” of the dataset SST, with algorithm PCA, and error parameter  $e = 30$ , which is highlighted in Figure 3.1 with a red circle.

The experimental results presented in this section suggest that if we were interested in compressing a dataset with many gaps, we would benefit from using the masking variant of an algorithm,  $a_M$ . However, even if the dataset didn’t have any gaps, the performance would not be significantly worse than that obtained by using the non-masking variant of the algorithm,  $a_{NM}$ . Therefore, since masking variants are, in general, more robust in this sense, in the sequel we focus on the set of variants  $V^*$  that we define next.

**Definition 3.2.2.** We denote by  $V^*$  the set of all the masking algorithm variants  $a_M$  for  $a \in A$ .

Notice that  $V^*$  includes a single variant for each algorithm. Therefore, in what follows we sometimes refer to the elements of  $V^*$  simply as algorithms.



### 3.3 Window Size Parameter

In this section, we analyze the extent to which the window size parameter impacts on the performance of the coding algorithms. For these experiments we consider the set of algorithm variants  $V_W^*$ , which is obtained from  $V^*$  by discarding algorithm SF, which doesn't have a window size parameter (recall this information from Table 2.1). Also, we only consider the four datasets that consist of multiple files, i.e. IRKIS, SST, ADCP and Solar (recall this information from Table 1.1). For each file, we compare the compression performance when using the OWS for the dataset, as defined in (3.6), and the LOWS for the file, defined next.

**Definition 3.3.1.** The *local optimal window size (LOWS)* of a coding algorithm variant  $a_v \in V_W^*$ , and a threshold parameter  $e \in E$ , for the data type  $z$  of a certain file  $f$  is given by

$$LOWS(a_v, e, z, f) = \arg \min_{w \in W} \left\{ CR(c_{\langle a_v, w, e \rangle}, z, f) \right\}, \quad (3.7)$$

where we break ties in favor of the smallest window size.

For each data type  $z$  of each dataset  $d$ , and each file  $f \in F(d, z)$ , coding algorithm variant  $a_v \in V_W^*$ , and threshold parameter  $e \in E$ , we calculate the RD between  $c_{\langle a_v, w_{global}^*, e \rangle}$  and  $c_{\langle a_v, w_{local}^*, e \rangle}$ , as defined in (3.2), where  $w_{global}^* = OWS(a_v, e, z, d)$  and  $w_{local}^* = LOWS(a_v, e, z, f)$ . In what follows, we denote the OWS and the LOWS as  $w_{global}^*$  and  $w_{local}^*$ , respectively.

As an example, in figures 3.3 and 3.4 we show  $w_{global}^*$ ,  $w_{local}^*$ , and the RD between  $c_{\langle a_v, w_{global}^*, e \rangle}$  and  $c_{\langle a_v, w_{local}^*, e \rangle}$ , as a function of the threshold parameter  $e$ , obtained for the data type  $z = \text{"VWC"}$ , for two different files of the dataset  $d = \text{IRKIS}$ . Figure 3.3 shows the results for the file  $f = \text{"vwc\_1202.dat.csv"}$ , and Figure 3.4 shows the results for  $f = \text{"vwc\_1203.dat.csv"}$ . Observe that the values of  $w_{global}^*$  are the same for both figures, which is expected, since both are obtained from the same data type of the same dataset.

In Figure 3.3 we notice, for instance, that for algorithm APCA the OWS and LOWS values match for every threshold parameter  $e$ , except 3 and 10. The OWS is larger than the LOWS when  $e = 3$ , but it is smaller when  $e = 10$ . In these two cases, the RD values are 1.52% and 1.76%, respectively. In Figure 3.4 we observe that in every case the OWS is larger than or equal to the LOWS. We highlight the marker for the maximum RD value (10.68%) obtained for all the tested CAIs, and we further comment on this point in the remaining of the section. Notice that in both figures the RD is non-negative in every plot, which makes sense, since the CR obtained with the OWS can never be lower than the CR obtained with the LOWS.

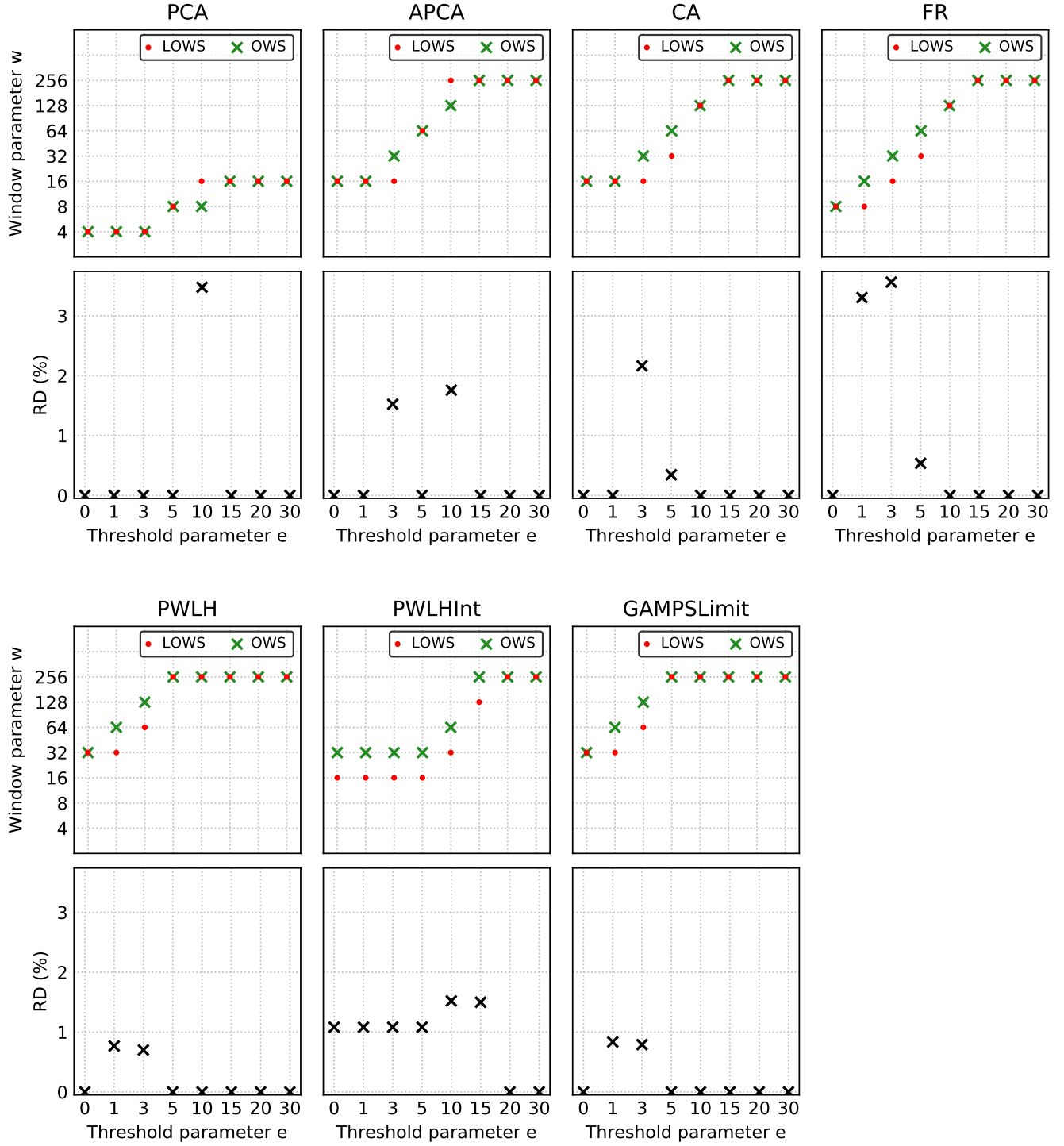


FIGURE 3.3: Plots of  $w_{global}^*$ ,  $w_{local}^*$ , and the RD between  $c_{<a_v, w_{global}^*, e>}$  and  $c_{<a_v, w_{local}^*, e>}$ , as a function of the threshold parameter  $e$ , obtained for the data type “VWC” of the file “vwc\_1202.dat.csv” of the dataset IRKIS.

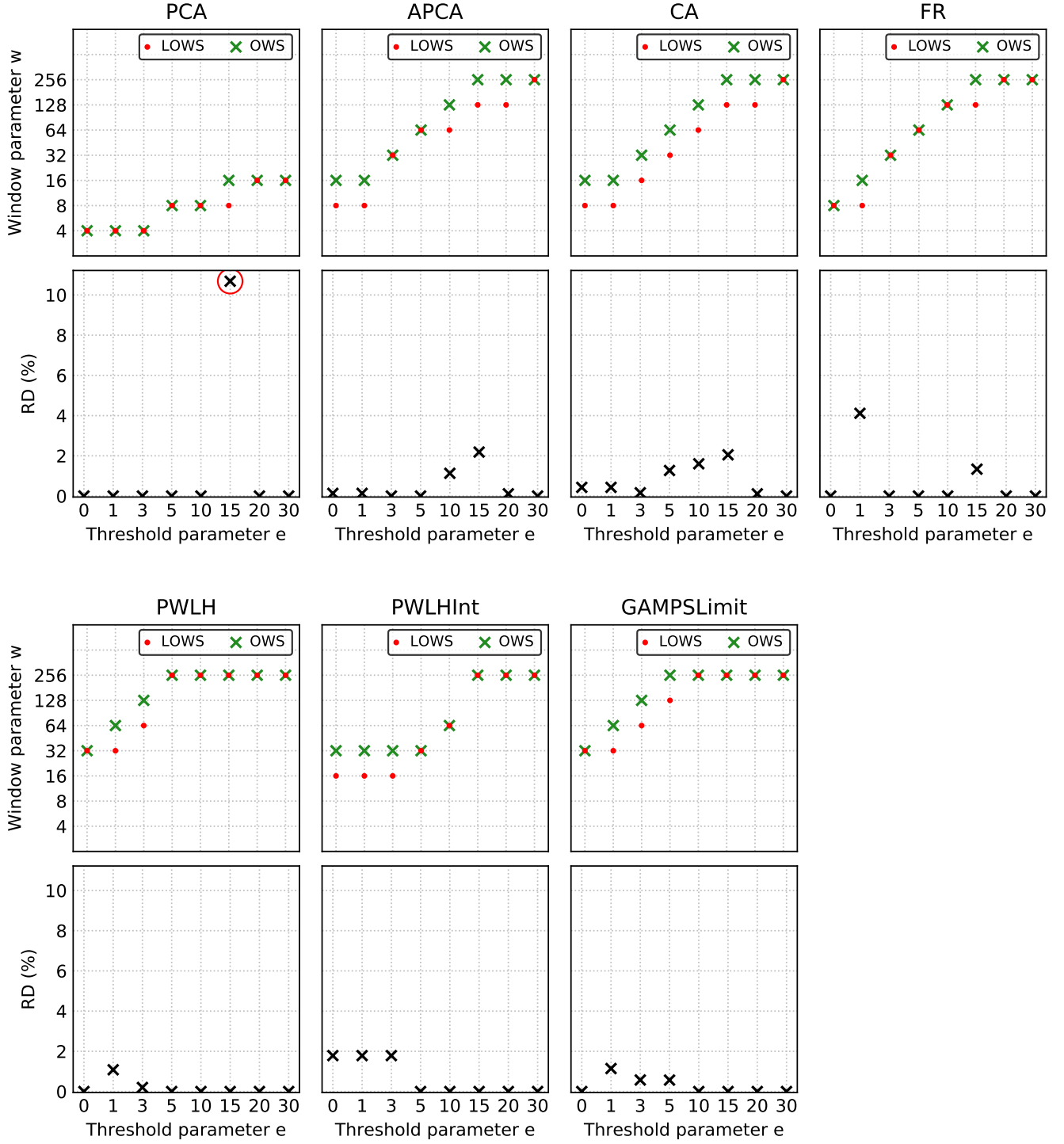


FIGURE 3.4: Plots of  $w_{global}^*$ ,  $w_{local}^*$ , and the RD between  $c_{<a_v, w_{global}^*, e>}$  and  $c_{<a_v, w_{local}^*, e>}$ , as a function of the threshold parameter  $e$ , obtained for the data type “VWC” of the file “vwc\_1203.dat.csv” of the dataset IRKIS. In the RD plot for algorithm PCA we highlight with a red circle the marker for the maximum value (10.68%) obtained for all the tested CAIs.

We analyze the experimental results to evaluate the impact of using the OWS instead of the LOWS on the compression performance of the tested coding algorithms. For each algorithm, we iterate through each threshold parameter, and each data type of each file, and we calculate the RD between the CAI with the OWS and the CAI with the LOWS. Since we consider 8 threshold parameters and there are 13 files with a single data type and 4 files with 3 different data types each, for each algorithm we compare a total of  $8 \times (13 + 4 \times 3) = 200$  pairs of CAIs. Table 3.2 summarizes the results of these comparisons, aggregated by algorithm and the range to which the RD belongs.

Algorithm	RD (%) Range				
	0	(0,1]	(1,2]	(2,5]	(5,11]
PCA	186 (93%)	4 (2%)	3 (1.5%)	2 (1%)	5 (2.5%)
APCA	174 (87%)	13 (6.5%)	7 (3.5%)	6 (3%)	0
CA	172 (86%)	16 (8%)	6 (3%)	6 (3%)	0
FR	171 (85.5%)	14 (7%)	8 (4%)	7 (3.5%)	0
PWLH	184 (92%)	13 (6.5%)	3 (1.5%)	0	0
PWLHInt	173 (86.5%)	9 (4.5%)	13 (6.5%)	4 (2%)	1 (0.5%)
GAMPSLimit	182 (91%)	16 (8%)	2 (1%)	0	0
Total	1,242 (88.7%)	85 (6.1%)	42 (3%)	25 (1.8%)	6 (0.4%)

TABLE 3.2: RD between the OWS and LOWS variants of each CAI.  
The results are aggregated by algorithm and the range to which the RD belongs.

For example, consider the results for algorithm CA, in the third row. The first column indicates that the RD is equal to 0 for exactly 172 (86%) of the 200 evaluated pairs of CAIs for that algorithm. The second column reveals that for 16 pairs of CAIs (8%), the RD takes values greater than 0 and less than or equal to 1%. The remaining three columns cover other ranges of RD values. Notice that for every row (except the last one), the values add up to a total of 200, since we compare exactly 200 pairs of CAIs for each algorithm.

The last row of Table 3.2 is obtained by adding the values of the previous rows, which combines the results for all algorithms. We notice that in 88.7% of the total number of evaluated pairs of CAIs, the RD is equal to 0. In these cases, in fact, the OWS and the LOWS coincide. In 97.8% of the cases, the RD is less than or equal to 2%. This means that, for the vast majority of CAI pairs, either the OWS and the LOWS match or they yield roughly the same compression performance. This result suggests that we could fix the window size parameter in advance, for example by optimizing over a training set, without compromising the performance of the coding algorithm. This is relevant, since calculating the LOWS for a file is, in general, computationally expensive.

We notice that there are only 6 cases (0.4%) in which the RD falls in the range (5, 11], most of which (5 cases) involve the algorithm PCA. The maximum value taken by RD (10.68%) is obtained for the data type “VWC” of the file “vwc\_1203.dat.csv” of the dataset SST, with algorithm PCA, and error parameter  $e = 15$ . In Figure 3.4 we highlight this maximum value with a red circle. In this case, the OWS is 16 and the LOWS is 8. According to these results, the performance of algorithm PCA seems to be more sensible to the window size parameter than the rest of the algorithms. Except for these few cases, we observe that, in general, the impact of using the OWS instead of the LOWS on the compression performance of coding algorithms is rather small. Therefore, in the following section, in which we compare the algorithms performance, we always use the OWS.

### 3.4 Algorithms Performance

In this section, we compare the compression performance of the coding algorithms presented in Chapter 2, by encoding the various datasets introduced in Chapter 1. We begin by comparing the algorithms among each other and later we compare them with gzip, a popular lossless compression algorithm. We analyze the performance of the algorithms on complete datasets (not individual files), so we always apply definitions 3.1.6–3.1.9. Following the results obtained in sections 3.2 and 3.3, we only consider the masking variants of the evaluated algorithms (i.e. set  $V^*$ ), and we always set the window size parameter to the OWS (recall Definition 3.2.1).

For each data type  $z$  of each dataset  $d$ , and each coding algorithm variant  $a_v \in V^*$  and threshold parameter  $e \in E$ , we calculate the CR of  $c_{\langle a_v, w_{global}^*, e \rangle}$ , as defined in (3.4), where  $w_{global}^* = \text{OWS}(a_v, e, z, d)$ . The following definition is useful for analyzing which CAI obtains the best compression result for a specific data type.

**Definition 3.4.1.** Let  $z$  be a data type of a certain dataset  $d$ , and let  $e \in E$  be a threshold parameter. We denote by  $c^b(z, d, e)$  the *best CAI* for  $z, d, e$ , and define it as the CAI that minimizes the CR among all the CAIs in CI, i.e.,

$$c^b(z, d, e) = \arg \min_{c \in CI} \left\{ CR(c_{\langle a_v, w_{global}^*, e \rangle}, z, d) \right\}. \quad (3.8)$$

When  $c^b(z, d, e) = c_{\langle a_v^b, w_{global}^{b*}, e \rangle}$ , we refer to  $a^b$  and  $w_{global}^{b*}$  as the *best coding algorithm* and the *best window size* for  $z, d, e$ , respectively.

Our experiments include a total of 21 data types, in 8 datasets. As an example, in Figure 3.5 we show the CR and the window size parameter  $w_{global}^*$ , as a function of the threshold parameter, obtained for each algorithm, for the data type “SST” of the dataset ElNino. For each threshold parameter  $e \in E$ , we use blue circles to highlight the markers for the minimum CR value and the best window parameter (in the respective plots corresponding to the best algorithm). For instance, for  $e = 0$ , the best CAI achieves a CR equal to 0.33 using algorithm PCA with a window size of 256. So in this case, algorithm PCA is the best coding algorithm, and 256 is the best window size. For the remaining seven values for the threshold parameter, the blue circles indicate that in every case the best algorithm is APCA, and the best window size ranges from 4 up to 32.

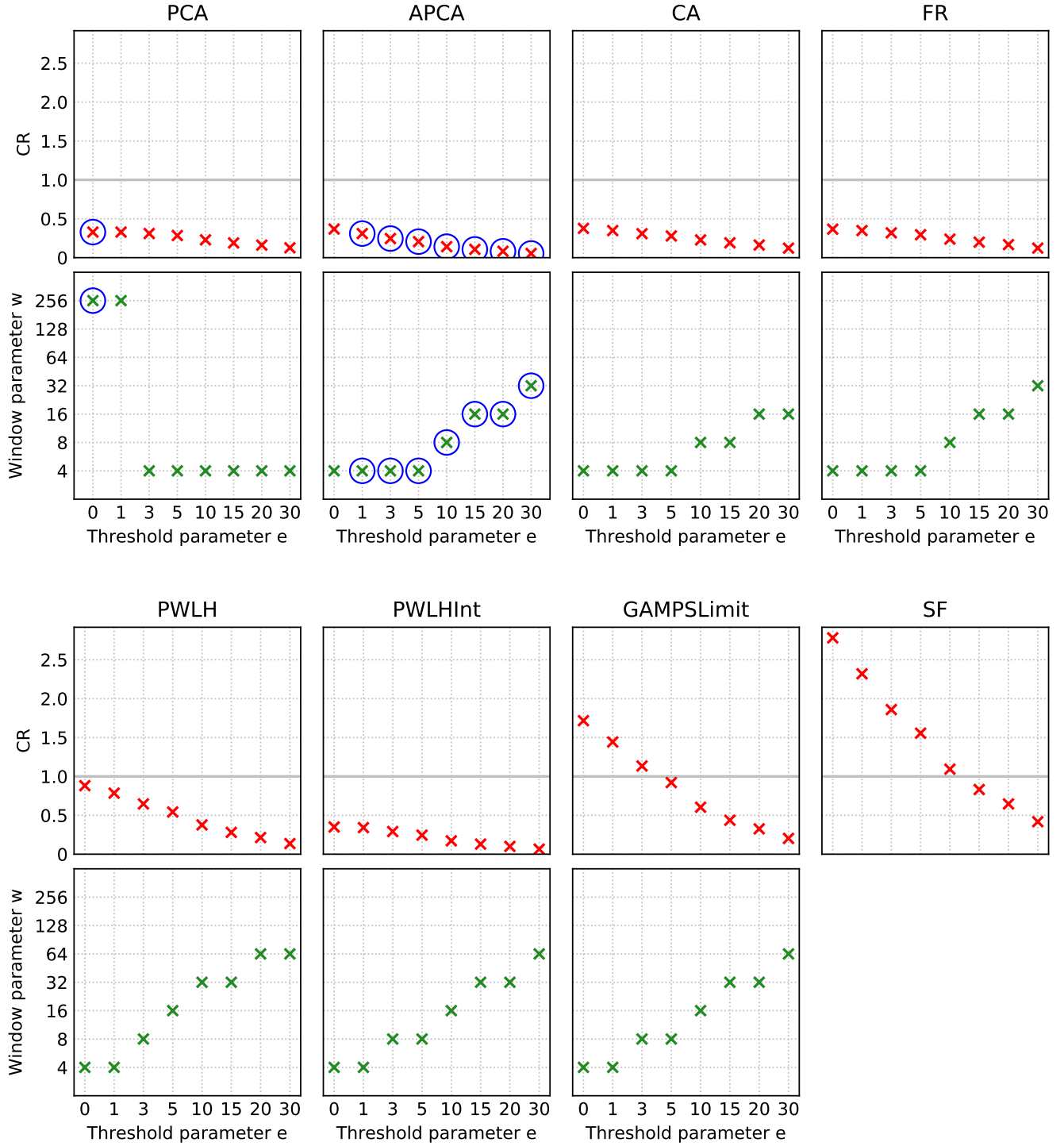


FIGURE 3.5: CR and window size parameter plots for every algorithm, for the data type “SST” of the dataset ElNino. For each threshold parameter  $e \in E$ , we use blue circles to highlight the markers for the minimum CR value and the best window size parameter (in the respective plots corresponding to the best algorithm)

Table 3.3 summarizes the compression performance results obtained by the evaluated coding algorithms, for each data type of each dataset. Each row contains information relative to certain data type. For example, the 13th row shows summarized results for the data type “SST” of the dataset ElNino, which are presented in more detail in Figure 3.5. For each threshold, the first column shows the CR obtained by the best CAI, the second column shows the base-2 logarithm of its window size parameter, and the cell color identifies the best algorithm.

		PCA		APCA		FR											
Dataset	Data Type	e = 0		e = 1		e = 3		e = 5		e = 10		e = 15		e = 20		e = 30	
		CR	w	CR	w	CR	w	CR	w	CR	w	CR	w	CR	w	CR	w
IRKIS	VWC	0.20	4	0.18	4	0.12	5	0.07	6	0.03	7	0.02	8	0.02	8	0.01	8
SST	SST	0.61	8	0.28	3	0.14	5	0.09	6	0.05	7	0.03	8	0.02	8	0.02	8
ADCP	Vel	0.68	8	0.68	8	0.67	2	0.61	2	0.48	2	0.41	2	0.35	3	0.26	3
Solar	GHI	0.78	2	0.76	3	0.71	4	0.67	4	0.59	4	0.52	4	0.47	4	0.38	4
	DNI	0.76	2	0.72	4	0.66	4	0.61	4	0.54	4	0.49	4	0.43	4	0.36	4
	DHI	0.78	2	0.77	2	0.72	4	0.68	4	0.60	4	0.54	4	0.48	4	0.39	4
ElNino	Lat	0.16	4	0.16	4	0.16	4	0.15	4	0.12	4	0.10	5	0.09	5	0.06	6
	Long	0.17	3	0.17	4	0.13	4	0.12	5	0.09	6	0.07	6	0.05	7	0.02	8
	Z. Wind	0.31	8	0.31	8	0.31	8	0.31	8	0.27	2	0.24	2	0.21	2	0.16	3
	M. Wind	0.31	8	0.31	8	0.31	8	0.31	8	0.29	2	0.26	2	0.23	2	0.19	2
	Humidity	0.23	8	0.23	8	0.23	8	0.23	8	0.21	2	0.18	2	0.16	2	0.13	2
	AirTemp	0.33	8	0.33	8	0.30	2	0.27	2	0.22	2	0.19	3	0.17	3	0.13	4
	SST	0.33	8	0.31	2	0.25	2	0.21	2	0.14	3	0.11	4	0.08	4	0.05	5
Hail	Lat	1.00	8	1.00	8	0.90	2	0.83	2	0.71	2	0.65	3	0.57	3	0.47	3
	Long	1.00	8	1.00	8	0.86	2	0.78	2	0.65	2	0.55	3	0.49	3	0.39	4
	Size	0.81	2	0.81	2	0.81	2	0.81	2	0.81	2	0.81	2	0.81	2	0.64	3
Tornado	Lat	1.00	8	0.85	2	0.71	2	0.65	2	0.54	3	0.47	3	0.42	4	0.33	4
	Long	1.00	8	0.82	2	0.65	2	0.58	3	0.46	3	0.40	4	0.35	4	0.28	4
Wind	Lat	1.00	8	1.00	8	0.89	2	0.81	2	0.70	2	0.62	3	0.56	3	0.47	3
	Long	1.00	8	0.95	2	0.80	2	0.73	2	0.62	3	0.54	3	0.49	3	0.40	4
	Speed	0.65	4	0.44	3	0.26	6	0.17	7	0.16	5	0.12	6	0.10	6	0.08	6

TABLE 3.3: Compression performance of the best evaluated coding algorithm, for various error thresholds on each data type of each dataset. Each row contains information relative to certain data type. For each threshold, the first column shows the minimum CR, and the second column shows the base-2 logarithm of the window size parameter for the best algorithm (the one that achieves the minimum CR), which is identified by a certain cell color described in the legend above the table.

We observe that there are only three algorithms (PCA, APCA, and FR) which are used by the best CAI for at least one of the 168 possible data type and threshold parameter combinations. Algorithm APCA is used in exactly 134 combinations (80%), including every case in which  $e \geq 10$ , and most of the cases in which  $e \in [1, 3, 5]$ . PCA is used in 31 combinations (18%), including most of the lossless cases, while FR is the best algorithm in only 3 combinations (2%), all of them for data type “Speed” of the dataset Wind.

Since there is not a single algorithm that obtains the best compression performance for every data type, it is useful to analyze how much is the RD between the best algorithm and the rest, for every experimental combination. With that in mind, next we define a pair of metrics.

**Definition 3.4.2.** The *maximum RD* ( $\text{maxRD}$ ) of a coding algorithm  $a \in A$  for certain threshold parameter  $e \in E$  is given by

$$\text{maxRD}(a, e) = \max_{z, d} \left\{ \text{RD}(c^b(z, d, e), c_{<a_v, w_{global}^*, e>}) \right\}, \quad (3.9)$$

where the maximum is taken over all the combinations of data type  $z$  and dataset  $d$ , and we recall that  $c^b(z, d, e)$  is the best CAI for  $z, d, e$ .

The maxRD metric is useful for assessing the compression performance of a coding algorithm  $a$  on the set of data types as a whole. Notice that maxRD is always non-negative. A satisfactory result (i.e. close to zero) can only be obtained when  $a$  achieves a good compression performance *for every data type*. In other words, bad compression performance *on a single data type* yields a poor result for the maxRD metric altogether. When maxRD is equal to zero,  $a$  achieves the best compression performance for every combination. Analyzing the results in Table 3.3, we observe that  $\text{maxRD}(\text{APCA}, e) = 0$  for every  $e \geq 10$ . Since the best algorithm is unique for every combination (i.e. exactly one algorithm obtains the minimum CR in every case), it is also true that, when  $a \neq \text{APCA}$ ,  $\text{maxRD}(a, e) > 0$  for every  $e \geq 10$ .

**Definition 3.4.3.** The *minmax RD* ( $\text{minmaxRD}$ ) for certain threshold parameter  $e \in E$  is given by

$$\text{minmaxRD}(e) = \min_{a \in A} \left\{ \text{maxRD}(a, e) \right\}, \quad (3.10)$$

and we refer to  $\arg \min_{a \in A}$  as the *minmax coding algorithm* for  $e$ .

Again, minmaxRD is always non-negative. Notice that  $\text{minmaxRD}(e) = 0$  for certain  $e$ , if and only if there exists a minmax coding algorithm  $a$  such that  $\text{maxRD}(a, e) = 0$ . Continuing the analysis from the previous paragraph, it should be clear that APCA is the minmax coding algorithm for every  $e \geq 10$ , since  $\text{maxRD}(\text{APCA}, e) = 0$  for every  $e \geq 10$ .

Table 3.4 shows the  $\text{maxRD}(a, e)$  obtained for every pair of coding algorithm variant  $a_v \in V^*$  and threshold parameter  $e \in E$ . For each  $e$ , the cell corresponding to the  $\text{minmaxRD}(e)$  value (i.e. the minimum value in the column) is highlighted.

Algorithm	maxRD (%)							
	e = 0	e = 1	e = 3	e = 5	e = 10	e = 15	e = 20	e = 30
PCA	40.52	42.28	53.11	62.01	71.73	75.33	77.21	80.28
APCA	33.25	15.64	9.00	29.96	0	0	0	0
CA	38.28	38.28	54.68	63.12	65.44	72.94	77.21	81.84
PWLH	73.46	72.93	72.52	82.14	83.24	86.86	88.94	91.19
PWLHInt	29.72	34.00	49.94	68.95	76.68	69.96	74.72	79.89
FR	48.75	49.85	52.21	52.70	54.82	55.35	54.48	64.72
SF	92.46	92.23	92.16	92.11	91.68	91.24	90.95	91.33
GAMPSLimit	85.84	85.73	84.37	83.92	83.02	83.00	82.88	82.22

TABLE 3.4:  $\text{maxRD}(a, e)$  obtained for every pair of coding algorithm variant  $a_v \in V^*$  and threshold parameter  $e \in E$ . For each  $e$ , the cell corresponding to the  $\text{minmaxRD}(a)$  value is highlighted.

In the lossless case, PWLHInt is the minmax coding algorithm, with minmaxRD being equal to 29.72%. This value is rather high, which means that none of the considered algorithms achieves a CR that is close to the minimum simultaneously *for every data type*. Recalling the results from Table 3.3 we notice that  $e = 0$  is the only threshold parameter value for which the minmax coding algorithm doesn't obtain the minimum CR in any combination. In other words, when



$e = 0$ , PWLHInt is the algorithm that minimizes the RD with the best algorithm among every data type, even though it itself is not the best algorithm for any data type.

When  $e \in [1, 3, 5]$ , the minmax coding algorithm is always APCA, and the minmaxRD values are 15.64%, 9.00% and 29.96%, respectively. Again, these values are fairly high, so we would select the most convenient algorithm depending on the data type we want to compress. Notice that in the closest case (algorithm FR for  $e = 5$ ), the second best maxRD (52.70%) is about 75% larger than the minmaxRD, which is a much bigger difference than in the lossless case.

When  $e \geq 10$ , the minmax coding algorithm is also always APCA, but in these cases the minmaxRD values are always 0. In the closest case (algorithm FR for  $e = 20$ ) the second best maxRD is 54.48%. If we wanted to compress any data type with any of these threshold parameter values, we would pick algorithm APCA, since according to our experimental results, it always obtains the best compression results with a significant difference over the remaining algorithms.

### 3.4.1 Comparison with algorithm gzip

In this subsection we consider the results obtained by the general purpose compression algorithm gzip [23]. This algorithm only operates in lossless mode (i.e. the threshold parameter can only be  $e = 0$ ), and it doesn't have a window size parameter  $w$ . Therefore, for each data type  $z$  of each dataset  $d$ , we have a unique CAI (and obtain a unique CR value) for gzip.

In all our experiments with gzip we perform a column-wise compression of the dataset files, which, in general, yields a much better performance than a row-wise compression. This is due to the fact that in most of our datasets, there is a greater degree of temporal than spatial correlation between the signals. All the reported results are obtained with the “--best” option of gzip, which targets compression performance optimization [24].

Table 3.5 summarizes the compression performance results obtained by gzip and the other evaluated coding algorithms, for each data type of each dataset. Similarly to Table 3.3, each row contains information relative to a certain data type, and for each threshold, the first column shows the CR obtained by the best CAI, the second column shows the base-2 logarithm of its window size parameter (when applicable), and the cell color identifies the best algorithm. Notice that, for  $e > 0$  we compare the gzip lossless result with the results obtained by lossy algorithms.

We observe that algorithm gzip obtains the best compression results in 36 (21%) of the 168 possible data type and threshold parameter combinations. Algorithms APCA, PCA, and FR now obtain the best results in exactly 106 (63%), 23 (14%), and 3 (2%) of the total combinations, respectively. Algorithm APCA is still the best algorithm for most of the cases in which  $e \geq 3$ . However, now there is no value of  $e$  for which APCA outperforms the rest of the algorithms for every data type, since gzip is the best algorithm for at least one data type in every case. In particular, gzip obtains the best compression results for the data type “Size” of the dataset Hail for every  $e$ . We also observe that gzip obtains the best relative results against the other algorithms for smaller values of  $e$ , which is expected, since lossy algorithms performance improves for larger values of  $e$ . However, even for  $e = 0$ , gzip only outperforms the rest of the algorithms in about a half (10 out of 21) of the data types.

		GZIP		PCA		APCA		FR									
Dataset	Data Type	e = 0		e = 1		e = 3		e = 5		e = 10		e = 15		e = 20		e = 30	
		CR	w	CR	w	CR	w	CR	w	CR	w	CR	w	CR	w	CR	w
IRKIS	VWC	0.13		0.13		0.12	5	0.07	6	0.03	7	0.02	8	0.02	8	0.01	8
SST	SST	0.52		0.28	3	0.14	5	0.09	6	0.05	7	0.03	8	0.02	8	0.02	8
ADCP	Vel	0.61		0.61		0.61		0.61	2	0.48	2	0.41	2	0.35	3	0.26	3
Solar	GHI	0.69		0.69		0.69		0.67	4	0.59	4	0.52	4	0.47	4	0.38	4
	DNI	0.67		0.67		0.66	4	0.61	4	0.54	4	0.49	4	0.43	4	0.36	4
	DHI	0.61		0.61		0.61		0.61		0.60	4	0.54	4	0.48	4	0.39	4
ElNino	Lat	0.08		0.08		0.08		0.08		0.08		0.08		0.08		0.06	6
	Long	0.07		0.07		0.07		0.07		0.07		0.07	6	0.05	7	0.02	8
	Z. Wind	0.31	8	0.31	8	0.31	8	0.31	8	0.27	2	0.24	2	0.21	2	0.16	3
	M. Wind	0.31	8	0.31	8	0.31	8	0.31	8	0.29	2	0.26	2	0.23	2	0.19	2
	Humidity	0.23	8	0.23	8	0.23	8	0.23	8	0.21	2	0.18	2	0.16	2	0.13	2
	AirTemp	0.33	8	0.33	8	0.30	2	0.27	2	0.22	2	0.19	3	0.17	3	0.13	4
	SST	0.32		0.31	2	0.25	2	0.21	2	0.14	3	0.11	4	0.08	4	0.05	5
Hail	Lat	1.00	8	1.00	8	0.90	2	0.83	2	0.71	2	0.65	3	0.57	3	0.47	3
	Long	1.00	8	1.00	8	0.86	2	0.78	2	0.65	2	0.55	3	0.49	3	0.39	4
	Size	0.37		0.37		0.37		0.37		0.37		0.37		0.37		0.37	
Tornado	Lat	1.00	8	0.85	2	0.71	2	0.65	2	0.54	3	0.47	3	0.42	4	0.33	4
	Long	1.00	8	0.82	2	0.65	2	0.58	3	0.46	3	0.40	4	0.35	4	0.28	4
Wind	Lat	1.00	8	1.00	8	0.89	2	0.81	2	0.70	2	0.62	3	0.56	3	0.47	3
	Long	1.00	8	0.95	2	0.80	2	0.73	2	0.62	3	0.54	3	0.49	3	0.40	4
	Speed	0.65	4	0.44	3	0.26	6	0.17	7	0.16	5	0.12	6	0.10	6	0.08	6

TABLE 3.5: Compression performance of the best evaluated coding algorithm, for various error thresholds on each data type of each dataset, including the results obtained by gzip. Each row contains information relative to certain data type. For each threshold, the first column shows the minimum CR, and the second column shows the base-2 logarithm of the window size parameter for the best algorithm (the one that achieves the minimum CR), which is identified by a certain cell color described in the legend above the table. Algorithm gzip doesn't have a window size parameter, so the cell is left blank in these cases.

Similarly to Table 3.4, Table 3.6 shows the  $\maxRD(a, e)$  obtained for every pair of coding algorithm variant  $a_v \in V^* \cup \{\text{gzip}\}$  and threshold parameter  $e \in E$ . For each  $e$ , the cell correspondent to the  $\min\maxRD(e)$  value is highlighted.

We observe that, for every  $e$ , the  $\min\maxRD$  values are rather high, the minimum being 26.58% (algorithm gzip for  $e = 0$ ). We conclude that none of the considered algorithms achieves a competitive CR *for every data type*, and the selection of the most convenient algorithm depends on the specific data type we are interested in compressing.

gzip is the  $\min\max$  coding algorithm when  $e \in [0, 1]$ , and in both cases the  $\min\maxRD$  values are rather high, i.e. 26.58% and 47.98%, respectively. APCA remains the  $\min\max$  coding algorithm for every  $e \geq 3$ , but its  $\min\maxRD$  values are now not only always greater than zero, but also quite high, ranging from 42.92% ( $e = 30$ ) up to 54.42% ( $e = 3$ ). This implies that there exist some data types for which the RD between the APCA and gzip CAIs is considerable, which means that, if we had the possibility of selecting gzip as a compression algorithm, APCA would no longer be the obvious choice for compressing any data type when  $e \geq 10$ , as we had concluded in the previous section.

Algorithm	maxRD (%)							
	e = 0	e = 1	e = 3	e = 5	e = 10	e = 15	e = 20	e = 30
GZIP	26.58	47.98	73.79	82.94	91.10	93.94	95.40	96.69
PCA	73.94	73.55	68.28	67.23	71.73	75.33	77.21	80.28
APCA	59.11	58.39	54.42	54.41	54.40	54.38	54.38	42.92
CA	73.83	73.46	69.31	68.30	65.44	72.94	77.21	81.84
PWLH	87.53	87.46	87.37	87.27	87.02	86.86	88.94	91.19
PWLHInt	71.26	71.01	70.71	68.95	76.68	69.96	74.72	79.89
FR	76.46	76.12	72.78	71.97	67.37	64.35	64.05	64.72
SF	96.31	96.13	96.09	95.96	95.87	95.74	95.64	95.05
GAMPSLimit	91.51	91.51	91.51	91.51	91.50	91.50	91.50	88.85

TABLE 3.6:  $\text{maxRD}(a, e)$  obtained for every pair of coding algorithm variant  $a_v \in V^* \cup \{\text{gzip}\}$  and threshold parameter  $e \in E$ . For each  $e$ , the cell corresponding to the  $\text{minmaxRD}(a)$  value is highlighted.

### 3.5 Conclusions

In conclusion, our experimental results indicate that none of the implemented coding algorithms obtains a satisfactory compression performance in every scenario. This means that selection of the best algorithm is heavily dependent on the data type to be compressed and the error threshold that is allowed. In addition, we have shown that, in some cases, even a general compression algorithm such as gzip can outperform our implemented algorithms. In general, according to our results, algorithms APCA and gzip achieve better compression results for larger error thresholds, while PCA, APCA, FR and gzip are preferred for lower thresholds. Therefore, if one wishes to compress certain data type, our recommended way for choosing the appropriate algorithm is to select the best algorithm for said data type according to Table 3.6.

In our research we have also compared the compression performance of the coding algorithms' masking and non-masking variants. The experimental results show that on datasets with few or no gaps both variants have a similar performance, while on datasets with many gaps the masking variant always performs better, sometimes achieving a significative difference. We concluded that the masking variant of a coding algorithm is preferred, since it is more robust and performs better in general.

In addition, we have studied the extent to which the window size parameter **impacts/affects** the compression performance of the coding algorithms. We analyzed the compression results obtained when using optimal global and local window sizes. The experimental results reveal that the **impact/effect** of using the optimal global window size instead of the optimal local window size for each file is rather small.

### 3.6 Future Work (TODO)

Some ideas:

- Consider non-linear models (e.g. Chebyshev Approximation)
- Consider new datasets
- Consider other metrics (e.g. RMSE)
- Investigate why certain algorithms perform better on certain data types

- 
- Create universal coder, with every algorithm as a subroutine

# Bibliography

- [1] EnviDat - IRKIS Soil moisture measurements Davos. <https://www.envidat.ch/dataset/soil-moisture-measurements-davos>, 2019. [Accessed 29 September 2020].
- [2] NOAA - TAO Data Download. [https://tao.ndbc.noaa.gov/tao/data\\_download/search\\_map.shtml](https://tao.ndbc.noaa.gov/tao/data_download/search_map.shtml), 2016. [Accessed 29 September 2020].
- [3] SolarAnywhere - Data. <https://data.solaranywhere.com/Data>, 2020. [Accessed 29 September 2020].
- [4] El Nino Data Set. <https://archive.ics.uci.edu/ml/datasets/El+Nino>, 1999. [Accessed 29 September 2020].
- [5] Kaggle - Storm Prediction Center. <https://www.kaggle.com/jtennis/spctornado>, 2016. [Accessed 29 September 2020].
- [6] N.Q.V. Hung, H. Jeung, and K. Aberer. An Evaluation of Model-Based Approaches to Sensor Data Compression. *IEEE Transactions on Knowledge and Data Engineering*, 25(11):2434–2447, 2013.
- [7] T. Bose, S. Bandyopadhyay, S. Kumar, A. Bhattacharyya, and A. Pal. Signal Characteristics on Sensor Data Compression in IoT - An Investigation. *2016 13th Annual IEEE International Conference on Sensing, Communication, and Networking (SECON)*, pages 1–6, 2016.
- [8] I. Lazaridis and S. Mehrotra. Capturing Sensor-Generated Time Series with Quality Guarantees. *Proc. ICDE*, pages 429–440, 2003.
- [9] E. Keogh, K. Chakrabarti, S. Mehrotra, and M. Pazzani. Locally Adaptive Dimensionality Reduction for Indexing Large Time Series Databases. *ACM Transactions on Database Systems*, 27(2):188–228, 2002.
- [10] C. Buragohain, N. Shrivastava, and S. Suri. Space Efficient Streaming Algorithms for the Maximum Error Histogram. *Proc. ICDE*, pages 1026–1035, 2007.
- [11] H. Elmeleegy, A.K. Elmagarmid, E. Cecchet, W.G. Aref, and W. Zwaenepoel. Online Piecewise Linear Approximation of Numerical Streams with Precision Guarantees. *Proc. VLDB Endowment*, 2(1):145–156, 2009.
- [12] G.E. Williams. Critical Aperture Convergence Filtering and Systems and Methods Thereof. *US Patent 7,076,402*, Jul. 11, 2006.
- [13] J.A.M. Heras and A. Donati. Fractal Resampling: time series archive lossy compression with guarantees. *PV 2013 Conference*, 2013.
- [14] S. Gandhi, S. Nath, S. Suri, and J. Liu. GAMPS: Compressing Multi Sensor Data by Grouping and Amplitude Scaling. *Proc. ACM SIGMOD International Conference on Management of Data*, pages 771–784, 2009.

- [15] J.J. Rissanen. Generalized Kraft Inequality and Arithmetic Coding. *IBM Journal of Research and Development*, 20(3):198–203, 1976.
- [16] Arithmetic Coding + Statistical Modeling = Data Compression. <https://marknelson.us/posts/1991/02/01/arithmetic-coding-statistical-modeling-data-compression.html>, 1991. [Accessed 29 September 2020].
- [17] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley-Interscience, 2nd edition, 2006.
- [18] R.E. Krichevsky and V.K. Trofimov. The performance of universal encoding. *IEEE Transactions on Information Theory*, 27(2):199–207, 1981.
- [19] I.H. Witten, R.M. Neal, and J.G. Cleary. Arithmetic Coding for Data Compression. *Communications of the ACM*, 30(6):520–540, 1987.
- [20] Data Compression With Arithmetic Coding. <https://marknelson.us/posts/2014/10/19/data-compression-with-arithmetic-coding.html>, 2014. [Accessed 29 September 2020].
- [21] Solomon W. Golomb. Run-length encodings (Corresp.). *IEEE Transactions on Information Theory*, 12(3), 1966.
- [22] R. Graham. An Efficient Algorithm for Determining the Convex Hull of a Finite Planar Set. *Information Processing Letters*, 1:132–133, 1972.
- [23] GNU Gzip. <https://www.gnu.org/software/gzip/>, 2018. [Accessed 29 September 2020].
- [24] GNU Gzip Manual - Invoking gzip. [https://www.gnu.org/software/gzip/manual/html\\_node/Invoking-gzip.html](https://www.gnu.org/software/gzip/manual/html_node/Invoking-gzip.html), 2018. [Accessed 29 September 2020].