

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

October 9, 2020

Abstract

TODO

Table of Contents

Abstract	ii
Table of Contents	iii
Index of Figures	v
Index of Tables	vii
1 Introduction	2
2 Datasets	3
2.1 Overview	4
2.2 Dataset IRKIS	5
2.3 Dataset ADCP	6
2.4 Dataset ElNino	6
2.5 Dataset Solar	6
2.6 Dataset Hail	6
2.7 Dataset Tornado	6
2.8 Dataset Wind	6
3 Algorithms	7
3.1 Introduction	8
3.2 Implementation details	9
3.2.1 Gap Encoding in the Masking Variants	10
3.3 Algorithm Base	10
3.4 Algorithm PCA	11
3.4.1 Example	13
3.4.2 Non-masking (<i>NM</i>) variant	15
3.5 Algorithm APCA	16
3.5.1 Example	17
3.5.2 Non-masking (<i>NM</i>) variant	19
3.6 Algorithms PWLH and PWLHInt	20
3.7 Algorithm SF	28
3.8 Algorithm CA	29
3.9 Algorithm FR	35
3.10 Algorithm GAMPS	40
3.11 Other	41
4 Experimental Results	42
4.1 Experimental Setting	43
4.2 Comparison of Masking and Non-Masking Variants	45
4.3 Window Size Parameter	49
4.4 Algorithms Performance	53

4.4.1	Comparison with algorithm gzip	57
4.5	Conclusions	59
4.6	Future Work (TODO)	59
Bibliography		62

List of Figures

3.1	Coding pseudocode for the Constant and Linear model algorithms.	9
3.2	Markov process diagram	10
3.3	Coding routine pseudocode for algorithm Base.	11
3.4	Coding routine pseudocode for variant M of algorithm PCA.	11
3.5	Decoding routine pseudocode for variant M of algorithm PCA.	12
3.6	13
3.7	13
3.8	14
3.9	14
3.10	Coding routine pseudocode for variant NM of algorithm PCA.	15
3.11	Coding routine pseudocode for variant M of algorithm APCA.	16
3.12	Decoding routine pseudocode for variant M of algorithm APCA.	17
3.13	17
3.14	18
3.15	18
3.16	PWLH.code_column_ M pseudocode.	20
3.17	PWLH.decode_column_ M pseudocode.	21
3.18	22
3.19	22
3.20	23
3.21	23
3.22	24
3.23	24
3.24	25
3.25	25
3.26	26
3.27	26
3.28	Example SF	28
3.29	CA.code_column_ M pseudocode.	29
3.30	30
3.31	30
3.32	31
3.33	31
3.34	32
3.35	32
3.36	33
3.37	33
3.38	34
3.39	FR.code_column_ M pseudocode.	35
3.40	FR.decode_column_ M pseudocode.	36
3.41	push_points_indexes pseudocode.	37
3.42	37

3.43	38
3.44	38
3.45	39
4.1	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.	46
4.2	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.	47
4.3	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.	50
4.4	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.	51
4.5	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)	54

Índice de tablas

2.1	Example of a dataset csv file with the format we defined.	4
2.2	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.	5
2.3	Statistics of dataset IRKIS. The gaps were ignored when calculating the median, mean and standard deviation of the sample values.	6
3.1	Characteristics of the evaluated coding algorithms. For each algorithm, the table shows whether it supports lossless and near-lossless compression (second and third columns), its model type (fourth column), whether the masking (M) and non-masking (NM) variants apply (fifth and sixth columns), and whether the algorithm depends on a window size parameter (w) (last column).	9
4.1	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.	48
4.2	RD between the OWS and LOWS variants of each CAI. The results are aggregated by algorithm and the range to which the RD belongs.	52
4.3	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.	55
4.4	$\maxRD(a, e)$ 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 $\min\maxRD(a)$ value is highlighted.	56
4.5	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.	58
4.6	$\maxRD(a, e)$ 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 $\min\maxRD(a)$ value is highlighted.	59

Cambios de la versión anterior (9/9/2020) a esta versión (1/10/2020)

- Agregué los índices al principio del pdf (incluyendo Table of Contents, para mejorar la navegación)

Chapter 2: Algorithms

- Hice las mejores sugeridas en las sections 3.1 (Introduction), 3.2 (Implementation details), 3.3 (Algorithm Base) y 3.4 (Algorithm PCA).
- Teniendo en cuenta las mejoras sugeridas para la Section 3.4 (Algorithm PCA), reescribí la Section 3.5 (Algorithm APCA). El Chapter 2 está pronto para leer hasta la Section 3.5 inclusive .
- Hice una Section 3.11 nueva, con algunos párrafos que no estoy seguro dónde ponerlos.
- En las gráficas de los ejemplos antes tenía dos leyendas: sample value y encoded value. Agregué decoded value. Creo que en los ejemplos queda claro qué significa cada cosa (especialmente en el ejemplo de PWLH - ver Figure 3.27, donde para cierto timestamp, un encoded value float no necesariamente coincide con el decoded value entero), pero quizás debería agregar una explicación en algún lado. Cuando eran dos leyendas había una explicación pero me sugeriste borrarla.
- Pregunta de estilo. Los comienzos de las secciones Example de PCA y APCA (3.4.1 y 3.5.1) son muy similares. ¿Está mal repetir el texto de forma tan similar o puedo comenzar el resto de secciones Example (de cada algoritmo) de la misma manera?
- En la Section 3.4 (Algorithm PCA), agregué la Subsection 3.4.2, en donde explico algunos detalles de la NM variant del algoritmo. Agregué el pseudocódigo de la rutina de codificación, pero la rutina de decodificación las expliqué solo con palabras, creo que es suficiente.
- En la Section 3.5 (Algorithm APCA), agregué la Subsection 3.5.2, pero en este caso los detalles de la NM variant los expliqué con palabras, sin agregar ningún pseudocódigo. Me parece que con tener los pseudocódigos de las dos rutinas de codificación (M y NM variants), aunque sea para un único algoritmo (PCA), basta para entender bien el tipo de diferencias que hay en el código de las dos variantes de cada algoritmo. Y creo que en el resto de algoritmos alcanza con tener el pseudocódigo de la rutina de codificación de la M variant y hacer algunos comentarios de la NM variant pero sin agregar pseudocódigo (igual que como hice con APCA).
- Las rutinas de decodificación para los algoritmos de modelo lineal (los que codifican los valores de una ventana con una función lineal) son bastante similares en todos los casos (PWLH, SF, CA y FR). En general decodifican un par de puntos (x_0, y_0) y (x_n, y_n) , y sustituyen x_i ($x_0 < x_i < x_n$) en la ecuación del segmento de recta (definido por los puntos) para obtener los valores de y_1, \dots, y_{n-1} . Creo que alcanza con tener el pseudocódigo de la rutina de decodificación para uno solo de estos algoritmos. Con el resto hay diferencias menores, que se pueden explicar con palabras.

Chapter 1

Introduction

Reutilizar lo que se pueda de las presentaciones del Pedeciba:

2018

<https://docs.google.com/presentation/d/1EtYbM5shn685DfP9qLd2E89LBBwyIevcM6oUQuE8RJA/edit>

<https://docs.google.com/document/d/1rBx11Ka9GhvohEkdwrMEuGeWtiitKOnbLrZ7yoOjjuE/edit>

2019

<https://docs.google.com/presentation/d/19glwhXE3IjQgQr-LBK5XV1lECWNcqBoEMdC9AILMSaE/edit>

<https://docs.google.com/document/d/1c8W0dungTl2JVxHxCv0cXXZ-uuMtzEaMVbQssSzw6M/edit>

Chapter 2

Datasets

In this chapter we present all the datasets that were compressed as part of our experimental work, which is described in Chapter 4. In Section 2.1 we show an example of a dataset csv file and analyze its information. We also 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 eight sections we describe each of the datasets, including some of their statistics.

2.1 Overview

Our project is focused on studying the compression of multichannel signals with irregular sampling and data gaps. Therefore, every one of the datasets presented in this chapter consists of signals with either one or both of these characteristics. In Chapter 4 we present our experimental results, which make use of these datasets to analyze the performance of every one of the compression algorithms presented in Chapter 3.

Since the datasets were gathered from multiple sources [1–6], they originally had different formats, so we decided to transform them into a consistent, homogeneous format. The format we defined is purposely generic so it can be easily adapted to represent any kind of dataset, aside from the ones considered in this project. In Figure 2.1 we present an example of a dataset csv file with said format. A csv file is simply a delimited text file that uses commas to separate values, but in general is more convenient to display its contents in a table. The first three rows have dataset metadata, the fourth row has the labels of the data columns, and the remaining rows consist of the actual data. The second cell in the first row has a key that uniquely identifies a dataset. In this example, the csv file corresponds to dataset ElNino, which we further examine in Section 2.4.

DATASET:	ElNino					
TIME UNIT:	hours					
FIRST TIMESTAMP:	1980-03-07 00:00:00					
Time Delta	Lat	Long	ZonWinds	MerWinds	AirTemp	SST
0	-2	-10946	-68	7	2614	2624
24	-2	-10946	-49	11	N	N
24	-2	-10946	-45	22	N	N
48	-1	-10946	-38	19	N	2431
24	-2	-10946	-42	15	2557	2319
48	-2	-10946	-44	3	2472	2364
24	-2	-10946	-32	1	N	2434

TABLE 2.1: Example of a dataset csv file with the format we defined.

The values in the first data column, which we label “Time Delta”, represent the timestamps associated with the data from the rest of the columns. We often refer to this column as the timestamp column. It is always the first column, and it must be present in every csv file, since we consider that in every dataset, each data value must be associated to a certain timestamp. In practice, this timestamp may represent the time at which the value was read, transmitted, stored, etc. The second cell from the second row, and the second cell from the third row, describe the time unit and the initial timestamp, respectively, for the timestamp column. Therefore, in the csv presented in Figure 2.1, the first four timestamps are “1980-03-07 00:00:00”, “1980-03-08 00:00:00”, “1980-03-09 00:00:00” and “1980-03-11 00:00:00”. Notice that the difference between subsequent timestamps is not constant, varying between 24 and 48 hours, which means that the signals in this dataset have irregular sampling (recall that this is one of the characteristics that we are looking for in signals).

Besides the timestamp column, the csv presented in Figure 2.1 consists of six additional data columns, each representing a different data type. The first two columns represent the latitude and longitude coordinates, respectively, of a buoy floating in the ocean, while the last four columns represent readings of different physical magnitudes (i.e. wind velocity, air temperature and sea temperature), made by sensors set in said buoy. In general, the timestamp column can only consist of integer values, while the rest of the columns can have both integer values and the

character “N”. An integer value represents an actual data sample, whose range depends on the range and accuracy of the sampling instrument used for acquiring and storing the data, while character “N” represents a gap in the data. In practice, this data gap may occur when there’s an error acquiring, transmitting and/or storing the data. Notice that, in the example csv, there are some gaps in the last two data columns (recall that this another one of the characteristics that we are looking for in signals).

In Figure 2.2 we show relevant information regarding each of the eight datasets presented in this chapter. The second column indicates the characteristic of each dataset, in terms of the amount of gaps. The third column shows the number of files in each dataset. In some cases, the full dataset consisted of more files, but we decided to narrow them down so that the experiments would have a reasonable runtime. The last two columns show the number of data types in each dataset, and their names, respectively. In many datasets there are multiple data columns for a single data type. Additional details of every dataset are presented in the remaining sections of this chapter.

Dataset	Dataset Characterstic	#Files	#Types	Data Types
IRKIS [1, 2]	Many gaps	7	1	VWC
SST [3]	Many gaps	3	1	SST
ADCP [3]	Many gaps	3	1	Vel
Solar [4]	Many gaps	4	3	GHI, DNI, DHI
ElNino [5]	Many gaps	1	7	Lat, Long, Zonal Winds, Merid. Winds, Humidity, Air Temp., SST
Hail [6]	No gaps	1	3	Lat, Long, Size
Tornado [6]	No gaps	1	2	Lat, Long
Wind [6]	No gaps	1	3	Lat, Long, Speed

TABLE 2.2: 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.

2.2 Dataset IRKIS

Dataset IRKIS [1] consists of soil moisture measurements from stations installed along the Dischma valley, in the municipality of Davos, Switzerland, in the period between October 2010 and October 2013. This information is particularly useful to researchers who study the role of soil moisture in relation with the snowpack runoff and catchment discharge in high alpine terrain [2].

The dataset was gathered in seven stations, labeled 1202, 1203, 1204, 1205, 222, 333, and SLF2, which can be located in the map presented in [1]. In each station, the soil moisture data were measured by sensors installed at different depths, i.e. 10, 30, 50, 80 and 120 cm, below the surface. There were two sensors installed per depth, labeled A and B in the dataset files. Thus, in each station, 10 different data samples were stored for each timestamp. The sensors measure the *Volumetric Water Content (VWC)*, which is equal to the ratio of water volume to soil volume. Therefore, higher VWC values indicate a more moist soil.

In Table 2.3 we present relevant statistics of dataset IRKIS. The data gathered in each station is stored in a separate csv file, and every row in the table contains statistics of a different file. The first three columns show the number of rows, columns and total entries (i.e. total rows times total columns), respectively. The fourth column specifies the number of gaps, and its percentage of the total entries. The last five columns show the minimum, maximum, median, mean, and standard deviation, of the sample values.

Station	#Rows	#Cols	#Entries	#Gaps (%)	Min	Max	Mdn	Mean	SD
1202	26,305	10	263,050	125,190 (47.6)	240	541	428	417.1	48.8
1203	26,305	10	263,050	200,051 (76.1)	165	385	249	257.7	40.1
1204	26,305	10	263,050	178,657 (67.9)	218	464	298	313.3	70.0
1205	26,305	10	263,050	224,287 (85.3)	272	600	315	342.1	63.9
222	26,304	10	263,040	56,007 (21.3)	128	562	426	384.2	119.7
333	26,088	10	260,880	37,520 (14.4)	38	451	327	291.5	81.1
SLF2	26,305	10	263,050	43,049 (16.4)	215	580	352	360.9	65.1

TABLE 2.3: Statistics of dataset IRKIS. The gaps were ignored when calculating the median, mean and standard deviation of the sample values.

2.3 Dataset ADCP

2.4 Dataset ElNino

2.5 Dataset Solar

2.6 Dataset Hail

2.7 Dataset Tornado

2.8 Dataset Wind

Chapter 3

Algorithms

In this chapter we present the coding algorithms implemented and evaluated in the project. In Section 3.1 we give an overview of the different algorithms, including their parameters and masking variants. In Section 3.2 we provide some implementation details, including the pseudocodes for the coding and decoding subroutines that are common to most algorithms. We also describe the implementation of the gap encoding in the masking variant, which applies the KT estimator and arithmetic coding. In Section 3.3 we present algorithm Base, which is a trivial algorithm that serves as a base ground for the compression performance comparison developed in Chapter 4. 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.

3.1 Introduction

There exists literature analyzing the performance of the state-of-the-art algorithms used for sensor data compression [7, 8]. 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 the datasets presented in Chapter 2 consist of signals that miss either one or both of these characteristics. We propose a number of variants of state-of-the-art algorithms, which are able to encode this type of signals.

We focus on algorithms that support near-lossless compression. Near-lossless compression guarantees a maximum point-by-point error between the decompressed and the original files. The error threshold can be specified by the user to the coding routine via the ϵ parameter. Observe that, when ϵ is equal to 0, compression is lossless, i.e. the decompressed and the original files are identical.

The algorithms follow a model-based compression approach that compresses signals by exploiting correlation among signal samples taken at close times (*temporal correlation*) and, in some cases, among samples from various signals (*spatial correlation*). In addition to efficient compression performance, they offer some data processing features, like inferring uncertain sensor readings, detecting outliers, indexing, etc. [7]. The model-based techniques are classified into different categories, depending on the model type: *Constant models* approximate signals by piecewise constant functions, *Linear models* use linear functions, and *Correlation models* simultaneously encode multiple signals exploiting temporal and spatial correlation. There also exist *Nonlinear models*, which approximate signals by complex nonlinear functions, but algorithms that follow this technique do not support near-lossless compression and yield poor compression results [7].

For most algorithms we propose 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 this chapter. We point out that the gaps in the decoded file must always match the gaps in the original file, regardless of the value of the error threshold parameter (ϵ), which is only considered when encoding sample values. In Section 4.2 we compare the compression performance of both variants, *M* and *NM*, for every algorithm that supports both.

Most of the algorithms support a window size parameter, denoted w , which defines the size of the blocks into which the data are divided for encoding. In algorithm PCA, parameter w defines a *fixed block size*, while in the rest of the algorithms it defines a *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.

In Table 3.1 we outline some characteristics of the evaluated algorithms and the proposed variants. For each algorithm, the second and third columns indicate whether it supports lossless and near-lossless compression, respectively, the fourth column shows its model type, the fifth and sixth columns indicate if the masking (*M*) and non-masking (*NM*) variants apply, respectively, and the last column specifies if the algorithm depends on a window size parameter (w). 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 near-lossless encoding.

Algorithm	Lossless	Near-lossless	Model	M	NM	w
Base	x		Constant		x	
PCA [9]	x	x	Constant	x	x	x
APCA [10]	x	x	Constant	x	x	x
PWLH [11]/ PWLHInt	x	x	Linear	x	x	x
SF [12]	x	x	Linear	x		
CA [13]	x	x	Linear	x	x	x
FR [14]	x	x	Linear	x		x
GAMPS [15]	x	x	Correlation	x	x	x

TABLE 3.1: Characteristics of the evaluated coding algorithms. For each algorithm, the table shows whether it supports lossless and near-lossless compression (second and third columns), its model type (fourth column), whether the masking (M) and non-masking (NM) variants apply (fifth and sixth columns), and whether the algorithm depends on a window size parameter (w) (last column).

3.2 Implementation details

Figure 3.1 shows a general encoding scheme used for every Constant and Linear model algorithm. The decoding scheme is symmetric. 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), algorithm GAMPS follows a different scheme, which we present in Section 3.10.

In Figure 3.1, the inputs for the coding routine are a csv data file in the format presented in Chapter 2, a key (v) that describes the algorithm variant (either M or NM), and the maximum error threshold (ϵ) and window size (w) parameters. The output is a binary file, which represents the input file encoded with a compression algorithm using the specified variant and parameters.

```

input :  $in$ : csv data file to be encoded
         $v$ : variant ( $M$  or  $NM$ )
         $\epsilon$ : maximum error threshold
         $w$ : window size
output:  $out$ : binary file with the encoding of  $in$ 
1 Create output file  $out$ 
2 Encode an algorithm identification key, and parameter  $w$  (if applies)
3 Encode the header of the input file
4 Encode the number of rows and columns in the input file
5 Encode the timestamps column using a lossless code
6 if  $v == M$  then
7   | Encode gap locations
8 end
9 Encode each signal column of the input file separately, using a coding routine for a
  specific coding algorithm

```

FIGURE 3.1: Coding pseudocode for the Constant and Linear model algorithms.

The timestamps column, which is comprised of integers, is the first column in every csv data file, and it is also the first column to be encoded (line 5). This is done using a lossless code in which every integer is encoded independently, using a fixed number of bits. We focus on the

compression of the sample columns (i.e. the rest of the columns in the data file), and do not delve into the optimization of timestamp compression, which we leave for future work. When the masking variant of the algorithm is executed, the positions of the gaps in every data column are encoded, in line 7; the details are explained in Subsection 3.2.1.

3.2.1 Gap Encoding in the Masking Variants

We recall that the masking variant of an algorithm starts by losslessly encoding the position of all the gaps in the data (line 7 in Figure 3.1). We describe the position of the gaps by encoding a sequence of binary symbols, $x_1 \dots x_n$, each symbol indicating the presence (0) or absence (1) of a sample in a specific timestamp. To this end we use an arithmetic coder (AC) [16–18], which, provided with a sequential probability assignment $p(x_i | x_1 \dots x_{i-1})$ for each symbol x_i given the past symbols $x_1 \dots x_{i-1}$, $1 \leq i \leq n$, generates an encoding bit stream of length $-\log P(x_1 \dots x_n) + O(1)$, where $P(x_1 \dots x_n) = \prod_{i=1}^n p(x_i | x_1 \dots x_{i-1})$. This code length is optimal up to an additive constant.

As we recall from Chapter 2, 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 model the sequence of gaps as a first-order Markov process, with Krichevsky–Trofimov [19], which we define next for binary alphabets.

Definition 3.2.1. Given a per-state sequential probability assignment string x over an alphabet $A = \{0, 1\}$, the *Krichevsky–Trofimov (KT) probability assignment* assigns the following probabilities for each symbol position i , $1 \leq i \leq n$

$$p(0 | x_1 \dots x_{i-1}) = \frac{n_0 + 1/2}{i}, \quad p(1 | x_1 \dots x_{i-1}) = \frac{n_1 + 1/2}{i}, \quad (3.1)$$

where n_0 and n_1 denote the number of occurrences of 0 and 1 in $x_1 \dots x_{i-1}$, respectively.

A first-order Markov process has two states, S_0 and S_1 , and we say that x_i occurs in state S_b iff the previous symbol, x_{i-1} , equals b . We arbitrary let S_1 be the initial state (i.e. the state in which x_1 occurs). In Figure 3.2 we present a diagram for this Markov process. A KT probability assignment for a first order Markov process is obtained by applying (3.1) separately for the subsequence of symbols that occur in states S_0 and S_1 . This is implemented by maintaining two pairs of symbol occurrence counters, n_0 , n_1 , one pair for each state.

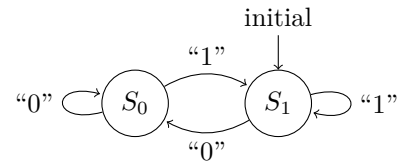


FIGURE 3.2: Markov process diagram

3.3 Algorithm Base

Algorithm Base is a trivial lossless coding 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 4.1.4), which we use to assess the compression performance of an algorithm in Chapter 4.

It follows the general schema presented in Figure 3.1, with a specific coding routine shown in Figure 3.3. This routine iterates through every column entry of a csv data file. Since algorithm Base only supports the *NM* variant, these entries can be either the character “N”, which represents a gap in the data, or an integer value representing an actual data sample. Every column

entry is encoded independently and using a fixed number of bits, which depends on the data type and the dataset. In practice, the number of bits used for encoding a data value ultimately depends on the range and accuracy of the sampling instrument used for acquiring and storing the data. A special integer, *NO_DATA*, is reserved for encoding a gap. The decoding routine is symmetric to the coding routine.

```

input : column: column of the csv data file to be encoded
         out: binary file encoded with algorithm Base
1 foreach entry in column, entry do
2   if entry == "N" then
3     value = NO_DATA
4   else
5     value = entry
6   end
7   Encode value using a (column-specific) fixed number of bits
8 end

```

FIGURE 3.3: Coding routine pseudocode for algorithm Base.

3.4 Algorithm PCA

Algorithm PCA [9], also known as Piecewise Constant Approximation, is a Constant model algorithm that supports lossless and near-lossless compression. It has a window size parameter (w) that establishes a fixed block size in which the data are separately processed and encoded. For PCA we define both variants, M and NM .

In Figure 3.4 we show the coding routine for variant M , in which all the column entries are integer values (there are no gaps). After adding the column entries into consecutive non-overlapping windows of size w (line 1), each of these windows is encoded independently (lines 3-13).

```

input : column: column of the csv data file to be encoded
         out: binary file encoded with algorithm PCA
          $\epsilon$ : maximum error threshold
          $w$ : fixed window size
1 Parse column into consecutive non-overlapping windows of size  $w$ , except possibly for
   the last window that may consist of fewer samples
2 foreach window in the parsing, win do
3   Let min and max be the minimum and maximum sample values in win, respectively
4   if  $|max - min| \leq 2 * \epsilon$  then
5     Output bit 0 to out
6     mid_range = (min + max)/2
7     Encode mid_range using a (column-specific) fixed number of bits
8   else
9     Output bit 1 to out
10    foreach value in win, value do
11      Encode value using a (column-specific) fixed number of bits
12    end
13  end
14 end

```

FIGURE 3.4: Coding routine pseudocode for variant M of algorithm PCA.

A window can be encoded in two different ways. If the absolute difference between its maximum and minimum values is less than or equal to $2 * \epsilon$ (i.e. the condition in line 4 is true), then bit 0 and the mid-range of the window are encoded (lines 5-7). On the other hand, if the condition in line 4 is false, then bit 1 and each of the window values are encoded independently (lines 9-12).

The decoding routine for variant M is shown in Figure 3.5. It keeps running until every entry in the column has been decoded, which occurs when condition in line 2 becomes false. Recall that the coding algorithm encodes the number of rows (line 4 in Figure 3.1), so this information is known by the decoding routine (as input *col_size*). First, a single bit is read from the input binary file (line 4). If that bit is 0, then the mid-range of an encoded window is decoded and written *size* times into the decoded csv data file (lines 6-7). On the other hand, if the bit read is 1, then the following process is repeated a total of *size* times: a fixed number of bits is read, and the associated value 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
        col_size: number of entries in the column

1  n = 0
2  while n < col_size do
3      size = min{w, col_size - n}
4      Decode bit from in
5      if bit == 0 then
6          Decode mid_range using a (column-specific) fixed number of bits
7          Output size copies of mid_range to out
8      else
9          repeat size times
10         Decode value using a (column-specific) fixed number of bits
11         Output value to out
12     end
13 end
14 n += size
15 end

```

FIGURE 3.5: Decoding routine pseudocode for variant M of algorithm PCA.

3.4.1 Example

Next we present an example of the encoding of 12 samples illustrated in Figure 3.6. Notice that the specific timestamp values are irrelevant for this algorithm. In this example we use algorithm PCA with an error threshold parameter (ϵ) equal to 1, and a fixed window size (w) equal to 4.

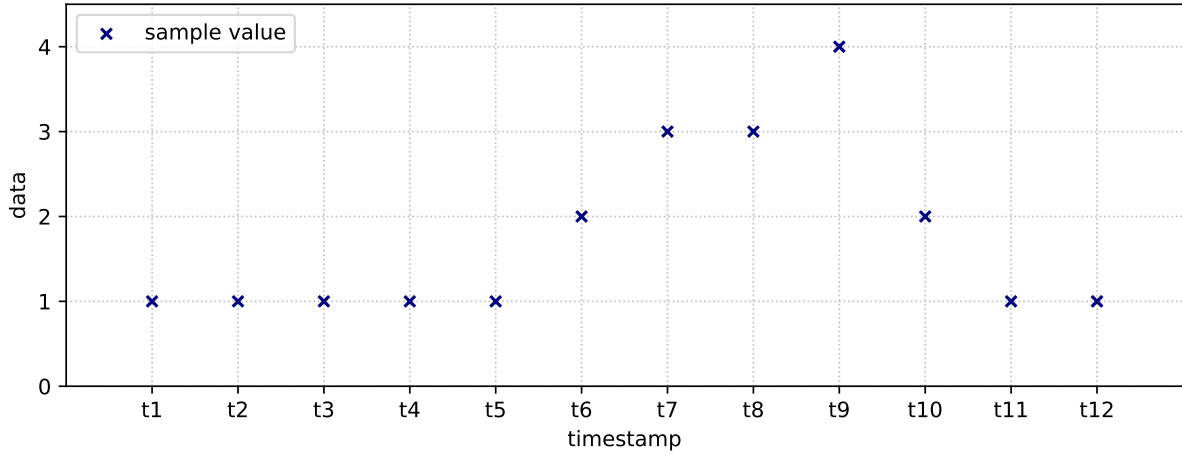


FIGURE 3.6

Since there are 12 samples to encode and $w = 4$, exactly three windows, each consisting of four samples, are encoded independently. The first window includes the first four samples, which are all equal to 1, so in this case the condition in line 4 of the coding routine is true, i.e. $|1 - 1| \leq 2 \cdot 1$. Therefore, the first window is encoded by executing lines 5-7. Since the mid-range of the window is 1, the first four values are encoded as 1. Figure 3.7 shows this step in the graph. Notice that, since all the values in the window are equal, the condition in line 4 would be true regardless of the value of parameter ϵ .

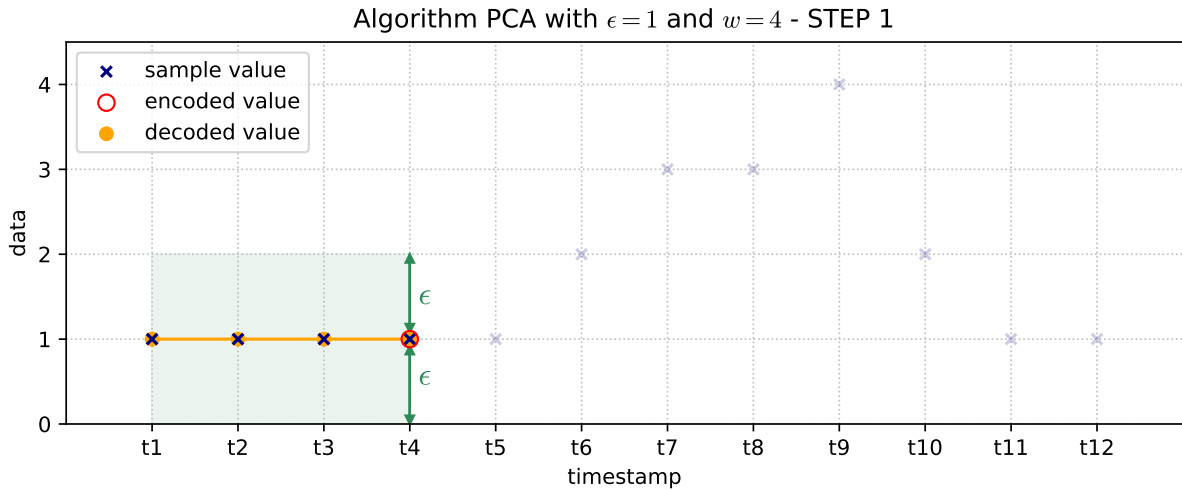


FIGURE 3.7

The second window is comprised of the next four samples, i.e. $[1, 2, 3, 3]$. Again, the condition in line 4 is true, since $|3 - 1| \leq 2 * 1$, but in this case the mid-range is 2, so these four values are encoded as 2. This step is shown in Figure 3.8.

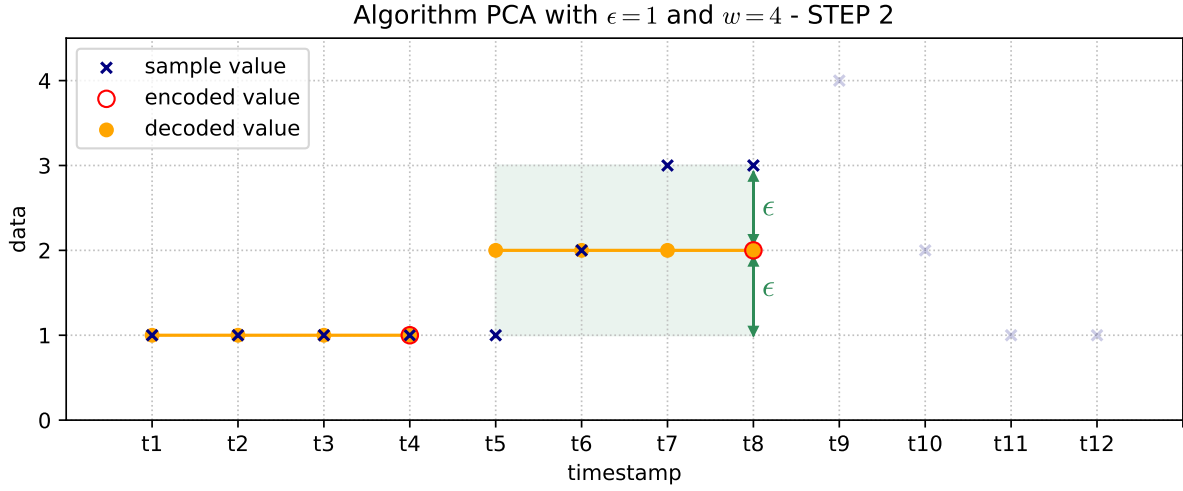


FIGURE 3.8

The third and last window consists of the last four samples, i.e. $[4, 2, 1, 1]$. In this case, the condition in line 4 is false, since $|4 - 1| > 2 * 1$, so the window is encoded by executing lines 9-12, which means that each of its values is encoded independently. This last step is shown in Figure 3.9.

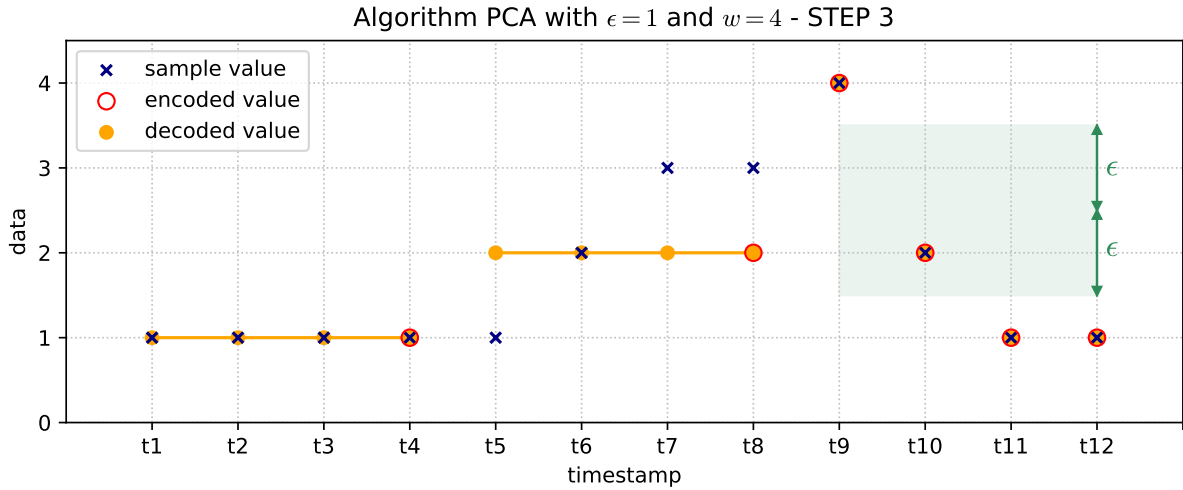


FIGURE 3.9

This simple example fairly represents every scenario that might arise during 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 illustrates why algorithm PCA is expected to achieve better compression performances on slowly varying signals rather than rough signals.

3.4.2 Non-masking (*NM*) variant

In Figure 3.10 we show the coding routine for variant *NM* of algorithm PCA. In this case, the column entries may be, not only an integer representing a sample value, but also the character “N” representing a gap in the data. As in variant *M*, after adding the column entries into consecutive non-overlapping windows of size w (line 1), each of these windows is encoded independently (lines 3-24). However, since not every entry in a window is guaranteed to be an integer, we must consider additional scenarios when encoding a window.

```

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

1 Parse column into consecutive non-overlapping windows of size  $w$ , except possibly for
  the last window that may consist of fewer samples
2 foreach window in the parsing, win do
3   if every entry in win is equal to “N” then
4     Output bit 0 to out
5     Encode NO_DATA using a (column-specific) fixed number of bits
6   else
7     Let min and max be the minimum and maximum sample values in win, resp.
8     if every entry in win is an integer and  $|max - min| \leq 2 * \epsilon$  then
9       Output bit 0 to out
10       $mid\_range = (min + max) / 2$ 
11      Encode mid_range using a (column-specific) fixed number of bits
12    else
13      Output bit 1 to out
14      foreach entry in win, entry do
15        if entry == “N” then
16          value = NO_DATA
17        else
18          value = entry
19        end
20        Encode value using a (column-specific) fixed number of bits
21      end
22    end
23  end
24 end

```

FIGURE 3.10: Coding routine pseudocode for variant *NM* of algorithm PCA.

A window can be encoded in three different ways. If every entry represents a gap in the data (i.e. the condition in line 3 is true), then bit 0 and the special integer *NO_DATA* are encoded (lines 4-5). If every entry in the window represents a sample value and the absolute difference between its maximum and minimum values is less than or equal to $2 * \epsilon$ (i.e. the condition in line 8 is true), then bit 0 and the mid-range of the window are encoded (lines 9-11). In every other case, bit 1 and each of the window entries are encoded independently (lines 13-21), using *NO_DATA* for encoding a gap. Notice that in the first two cases the window is encoded with the same amount of bits, i.e. $1 + column.total_bits$, while in the last case the window is encoded with $1 + w * column.total_bits$ bits.

The decoding routine for variant *NM* is quite similar to the decoding routine for variant *M*, which was presented in Figure 3.5, the only difference being that, in lines 6-7 and 10-11, when *NO_DATA* is decoded, it is character “N” what must be written into the decoded csv data file.

3.5 Algorithm APCA

Algorithm APCA [10], also known as Adaptive Piecewise Constant Approximation, is a Constant model algorithm that supports lossless and near-lossless compression. 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 separately processed and encoded is not fixed, but variable. In this case, the window size parameter (w) establishes the maximum block size allowed for the algorithm. APCA supports both variants, M and NM .

In Figure 3.11 we show the coding routine for variant M , in which all the column entries are integer values (there are no gaps). In every iteration an entry is added to the window (line 3), and a conditional that depends on parameters w and ϵ is checked (line 5). If the new entry makes the window violate the error threshold constraint (i.e. the absolute difference between its maximum and minimum values is greater than $2 * \epsilon$), or the window size restriction (i.e. its size is greater than w), then the entry is removed from the window (line 6), the window is encoded (lines 7-9), and a new window that includes the entry is created (lines 10). 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 column.total_bits bits are used for encoding its mid-range.

```

input : column: column of the csv data file to be encoded
         out: binary file encoded with algorithm APCA
          $\epsilon$ : maximum error threshold
          $w$ : maximum window size
1 Create a new window, win
2 foreach entry in column, entry do
3   Add entry to win
4   Let min and max be the minimum and maximum sample values in win, respectively
5   if  $|max - min| > 2 * \epsilon$  or win.size ==  $w + 1$  then
6     Remove entry from win, then recalculate min and max
7     Encode win.size using  $\log_2 w$  bits
8     mid_range =  $(min + max) / 2$ 
9     Encode mid_range using a (column-specific) fixed number of bits
10    Create a new window, win, then add entry to win
11  end
12 end

```

FIGURE 3.11: Coding routine pseudocode for variant M of algorithm APCA.

The decoding routine for variant M is shown in Figure 3.12. It keeps running until every entry in the column has been decoded, which occurs when condition in line 2 becomes false. The decoding loop is fairly simple. First, both the window size, $size$, and its mid-range value are decoded (lines 2-3). Then, the mid-range value is written $size$ times into the decoded csv data file (line 5).

```

input :  $in$ : binary file coded with algorithm APCA
          $out$ : decoded csv data file
          $w$ : maximum window size
          $col\_size$ : number of entries in the column
1  $n = 0$ 
2 while  $n < col\_size$  do
3   Decode  $size$  using  $\log_2 w$  bits
4   Decode  $mid\_range$  using a (column-specific) fixed number of bits
5   Output  $size$  copies of  $mid\_range$  to  $out$ 
6    $n += size$ 
7 end

```

FIGURE 3.12: Decoding routine pseudocode for variant M of algorithm APCA.

3.5.1 Example

Next we present an example of the encoding of 12 samples illustrated in Figure 3.6. Notice that the specific timestamp values are irrelevant for this algorithm. In this example we use algorithm APCA with an error threshold parameter (ϵ) equal to 1, and a maximum window size (w) equal to 256.

The condition in line 5 of the coding routine is false for the first eight iterations, so those samples, i.e. $[1, 1, 1, 1, 1, 2, 3, 3]$, are added to the first window. The sample processed in the 9th iteration is 4, which causes the window to violate the error threshold constraint, 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 3.13.

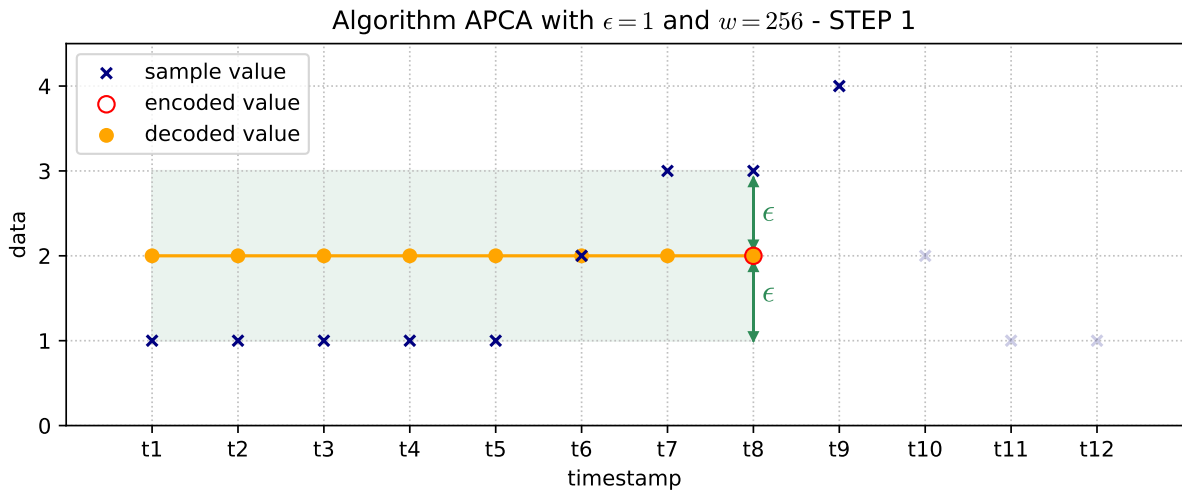


FIGURE 3.13

For the second window, the condition in line 5 is false in the 10th iteration. However, the error threshold constraint is violated in the 11th iteration, for the sample value 1, since again $|4 - 1| > 2 * 1$. That value is removed from the second window and 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 3.14.

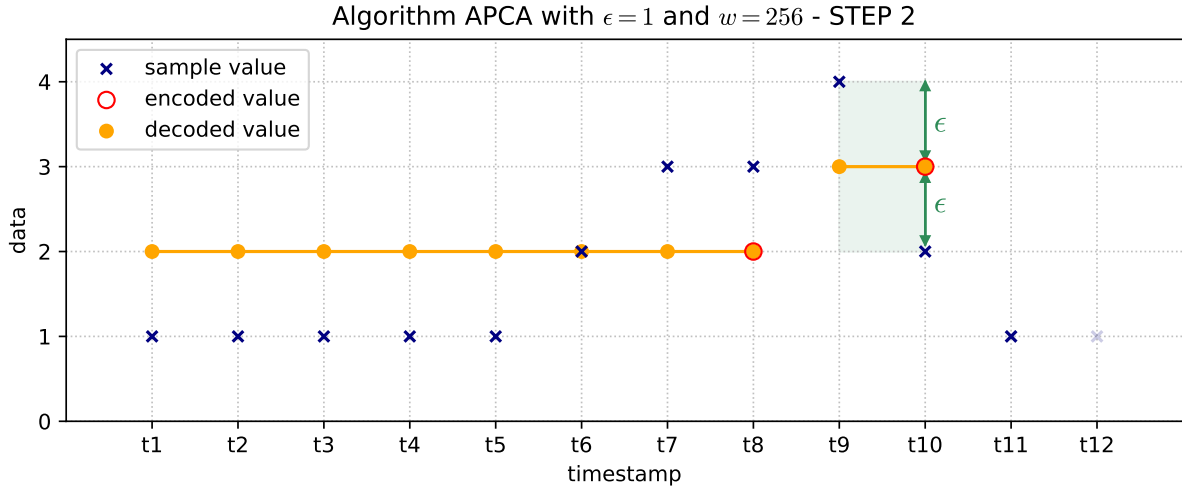


FIGURE 3.14

For the third window, the condition in line 5 is false in the 12th and last iteration. This means that said window, which is equal to $[1, 1]$, must be encoded after executing the last iteration. This is done by executing the same code as in lines 7-9, but was left out of the pseudocode for clarity. In this case, the window size is 2 and its mid-range is 1. This last step is shown in Figure 3.15.

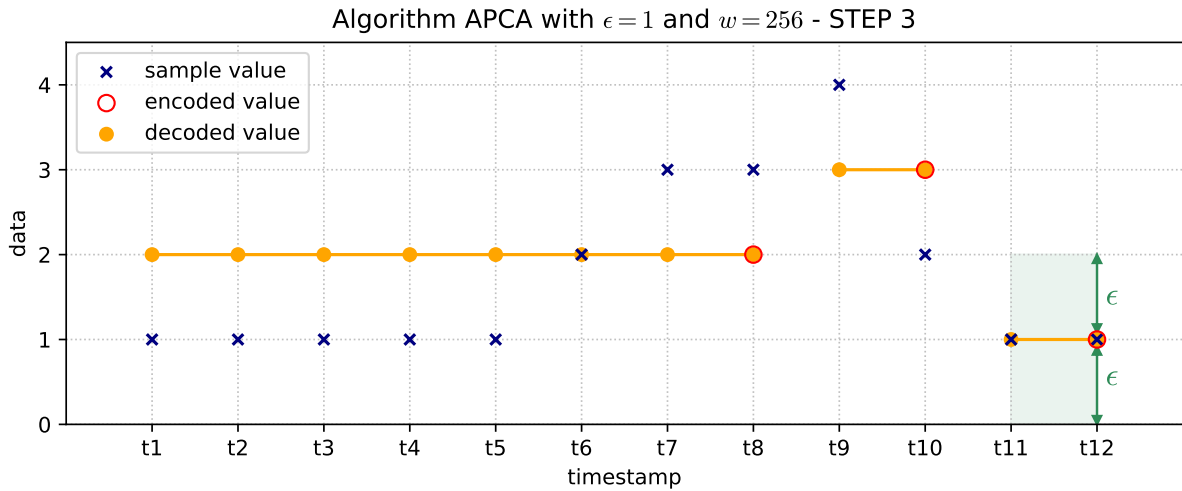


FIGURE 3.15

3.5.2 Non-masking (*NM*) variant

The coding and decoding routines for variant *NM* of algorithm APCA are similar to their variant *M* counterparts, the difference being that the former routines are able to handle both sample values *and gaps*. Recall that in the coding routine for variant *M*, a window is encoded when its newest entry makes it violate the error threshold constraint or the window size restriction (line 5 in Figure 3.11). In the coding routine for variant *NM*, a window must also be encoded if the newest entry is character “N” (gap in the data) and the other entries in the window are integers (sample values), or vice versa. A window that consists of gaps is encoded with the same amounts of bits as a window that consists of integers, i.e. $\log_2 w + \text{column.total_bits}$, where $\log_2 w$ bits are used for encoding its size, and column.total_bits bits are used for encoding the special integer *NO_DATA*.

3.6 Algorithms PWLH and PWLHInt

Algorithm PWLH [11], 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 [20], we reused the code from the framework linked in [7].

Figure 3.16 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 encoded
         out: binary file encoded with algorithm PWLH
          $\epsilon$ : maximum error threshold
          $w$ : maximum window size

1 Create a new window, win
2 foreach entry in column, entry 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      Remove entry from win
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    Create a new window, win
26    win.push(entry)
27    win.update_convex_hull()
28  end
29 end

```

FIGURE 3.16: 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 3.1), 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 3.17 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 3.17: 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 3.6, 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 (ϵ) 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 3.18 shows this step in the graph. Observe that, besides the original values, the convex hull for the current window is also shown.

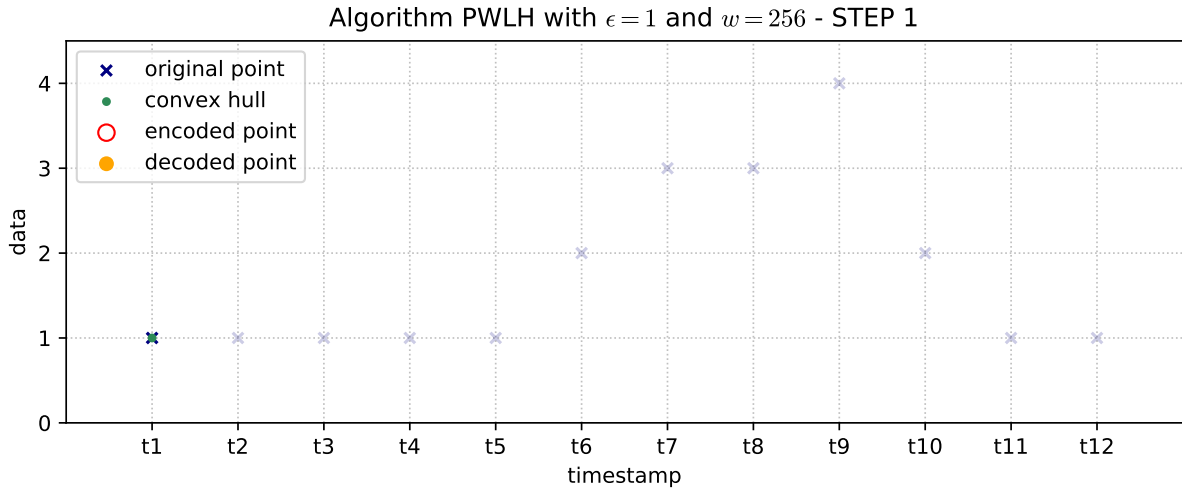


FIGURE 3.18

After adding the second data value to the window, the convex hull is updated. This step is shown in Figure 3.19. 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.

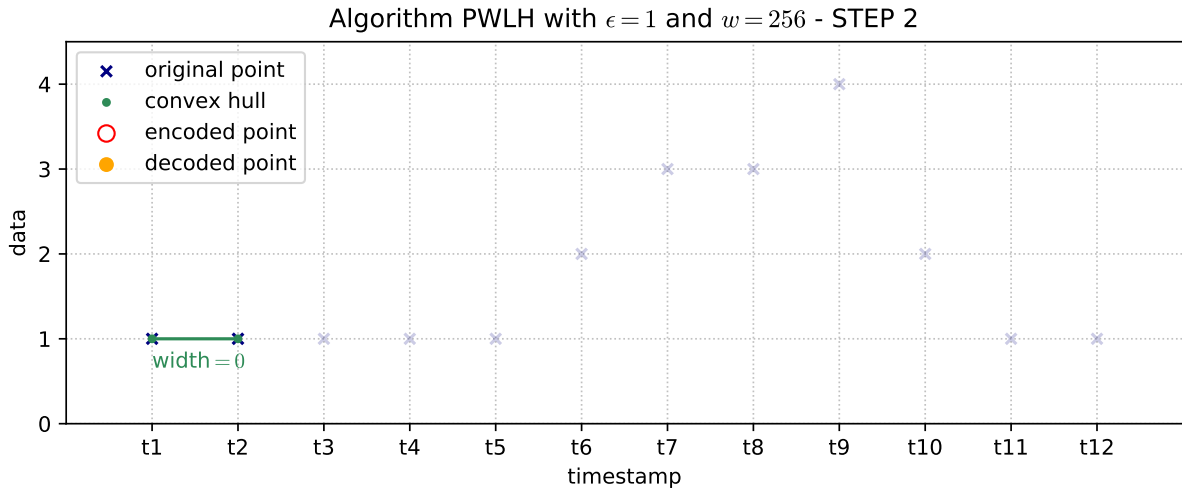


FIGURE 3.19

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 3.20.

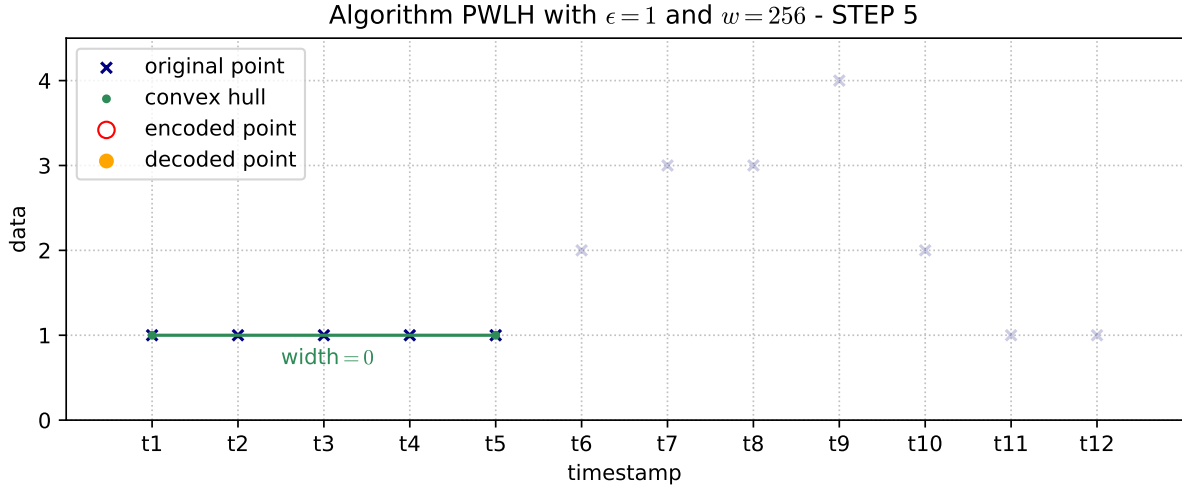


FIGURE 3.20

Next, the data value 2 is added to the window. The updated convex hull, which now consists of three edges, is shown in Figure 3.21. 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.

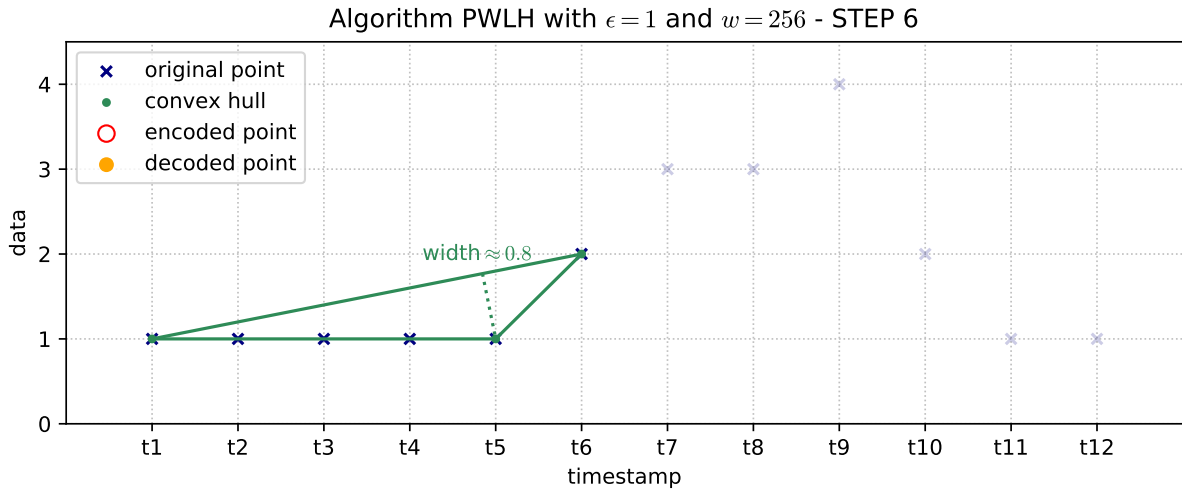


FIGURE 3.21

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 3.22, 3.23 and 3.24.

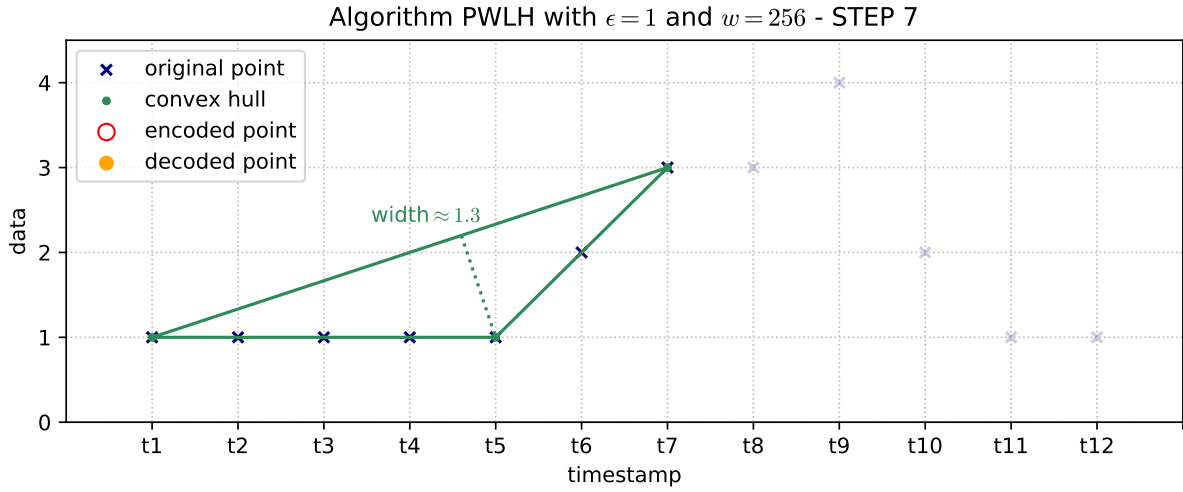


FIGURE 3.22

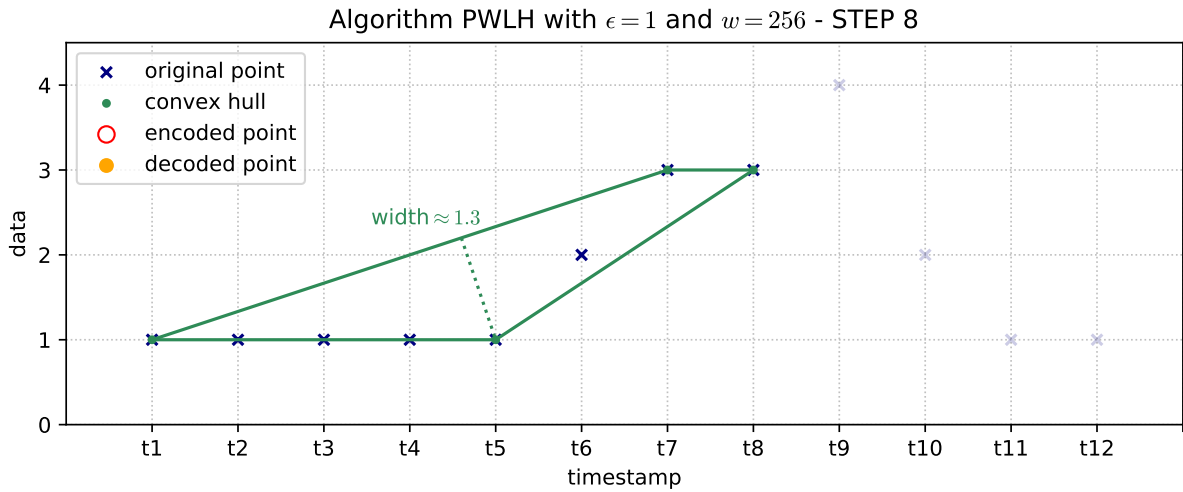


FIGURE 3.23

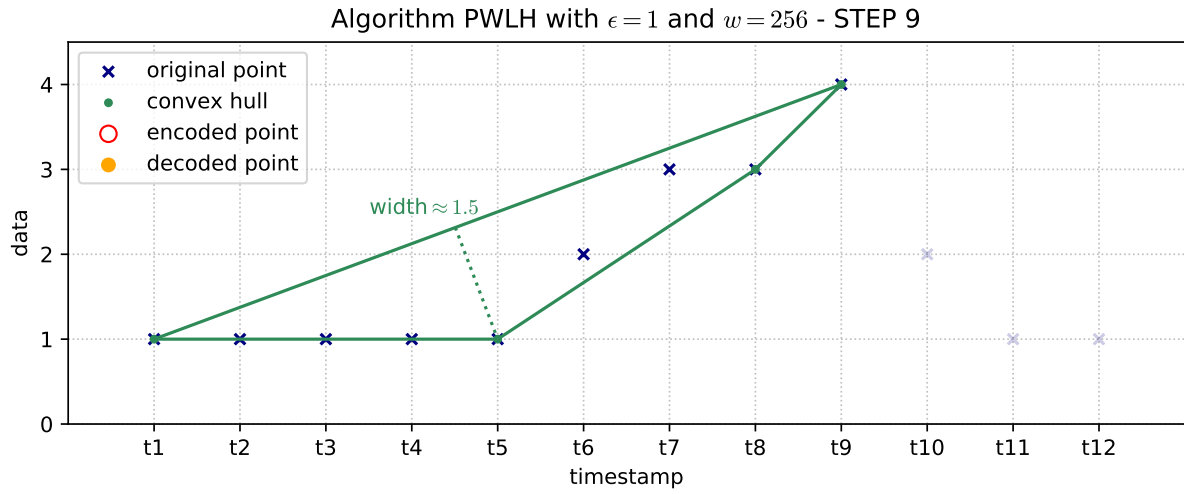


FIGURE 3.24

Eventually, in the 10th iteration, after adding the data value 2 to the window, the resulting convex hull, which is shown in Figure 3.25, 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.

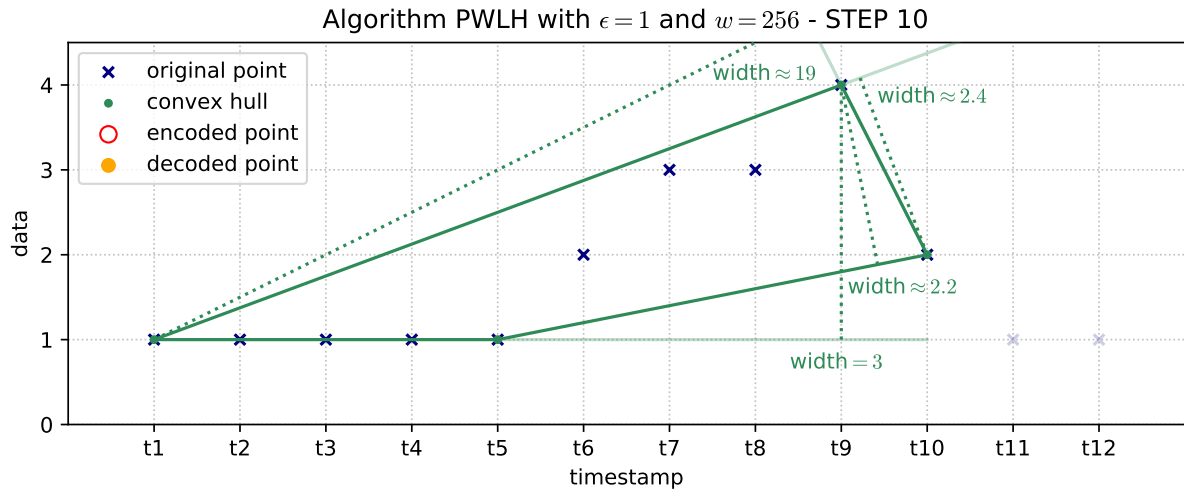


FIGURE 3.25

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 3.26 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.

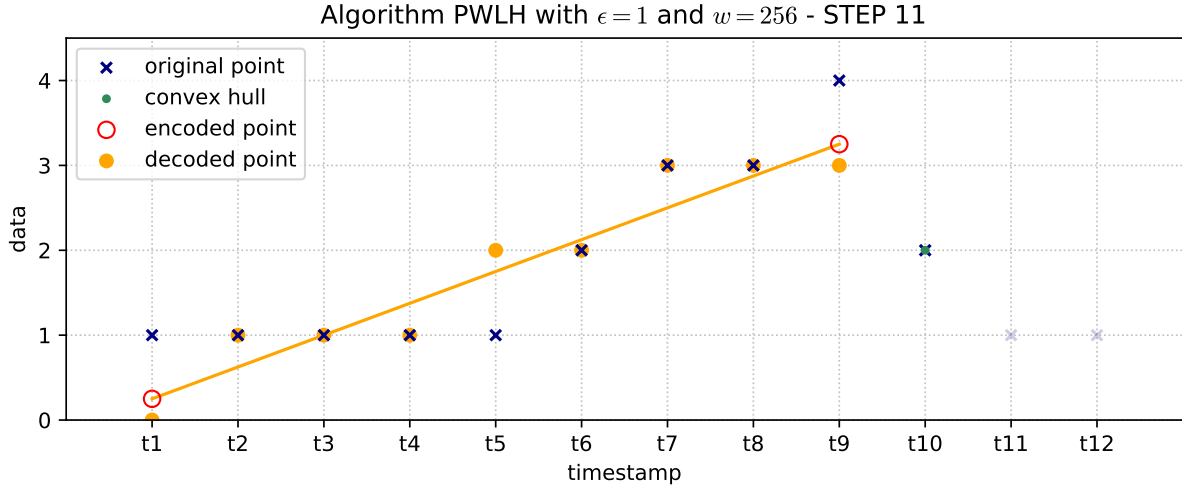


FIGURE 3.26

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 3.27 shows the line segment that minimizes the maximum distance to the three points in the second window.

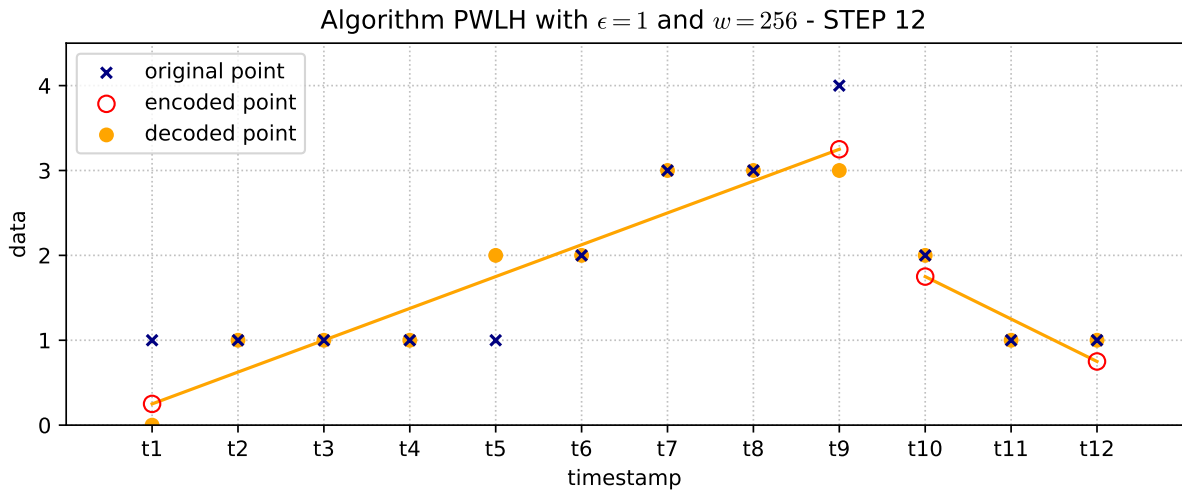


FIGURE 3.27

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

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

3.7 Algorithm SF

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

Details are omitted...

Describe example in Figure 3.28.

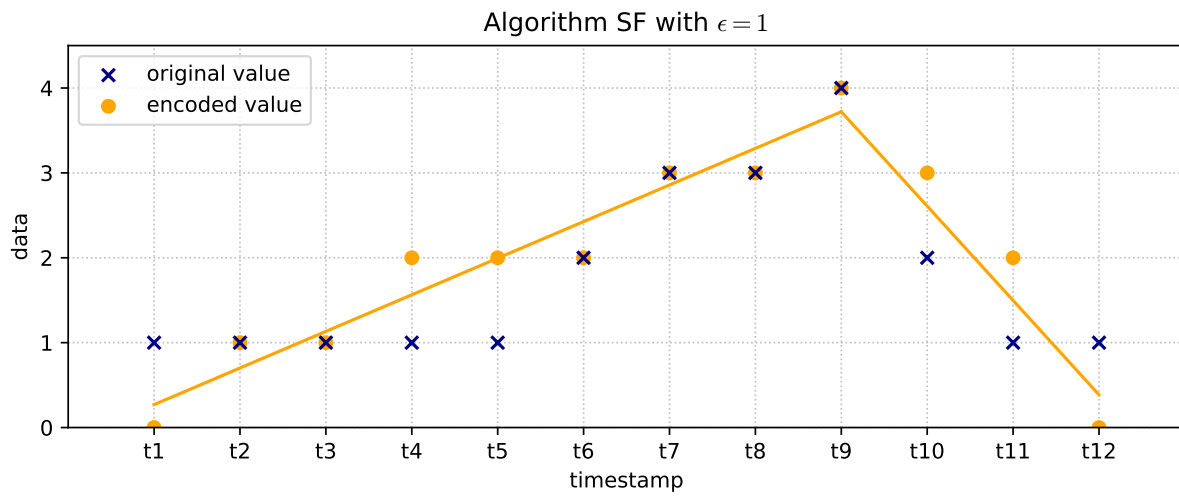


FIGURE 3.28: Example SF

3.8 Algorithm CA

Algorithm CA [13], 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 3.29 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 encoded
        out: binary file encoded with algorithm CA
         $\epsilon$ : maximum error threshold
         $w$ : maximum window size
1 Create a new window, win
2 foreach entry in column, entry 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 3.29: 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 $f(t_6) = f(t_7)$)

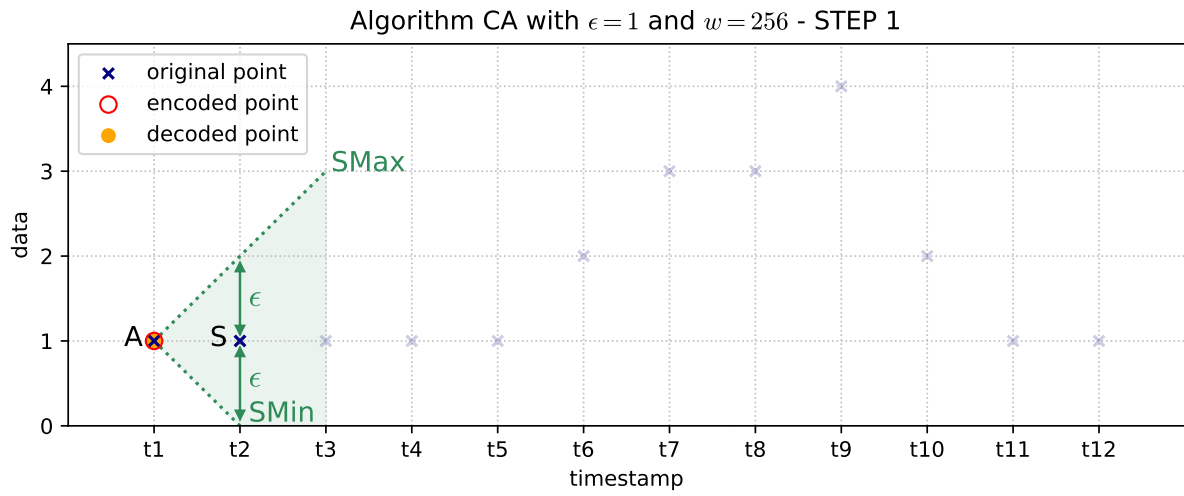


FIGURE 3.30

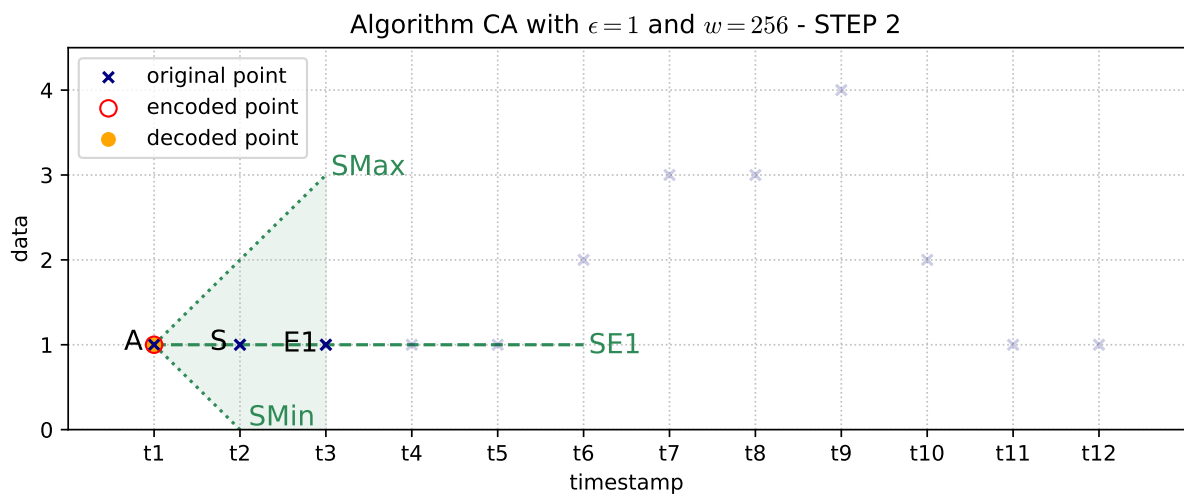


FIGURE 3.31

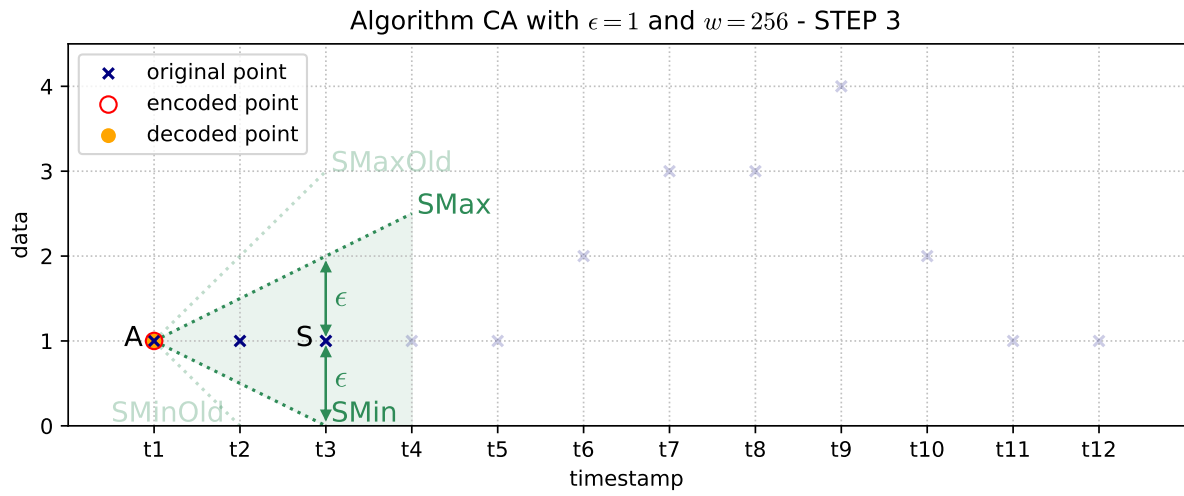


FIGURE 3.32

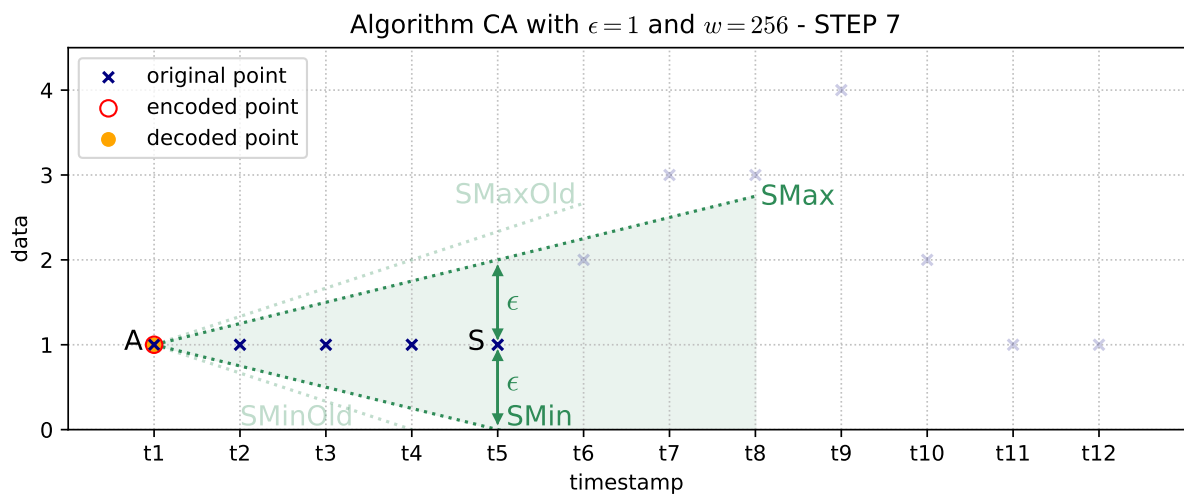


FIGURE 3.33

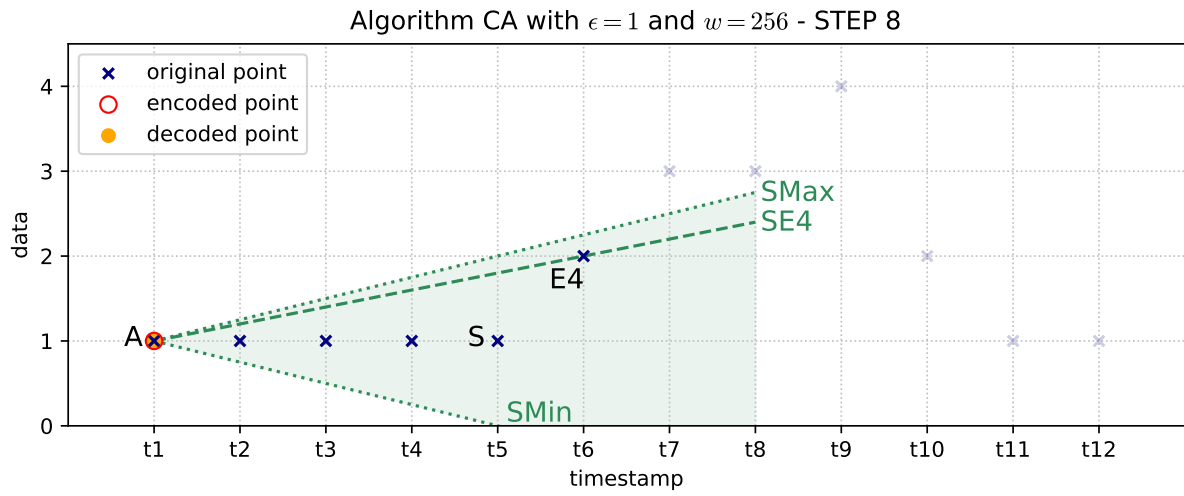


FIGURE 3.34

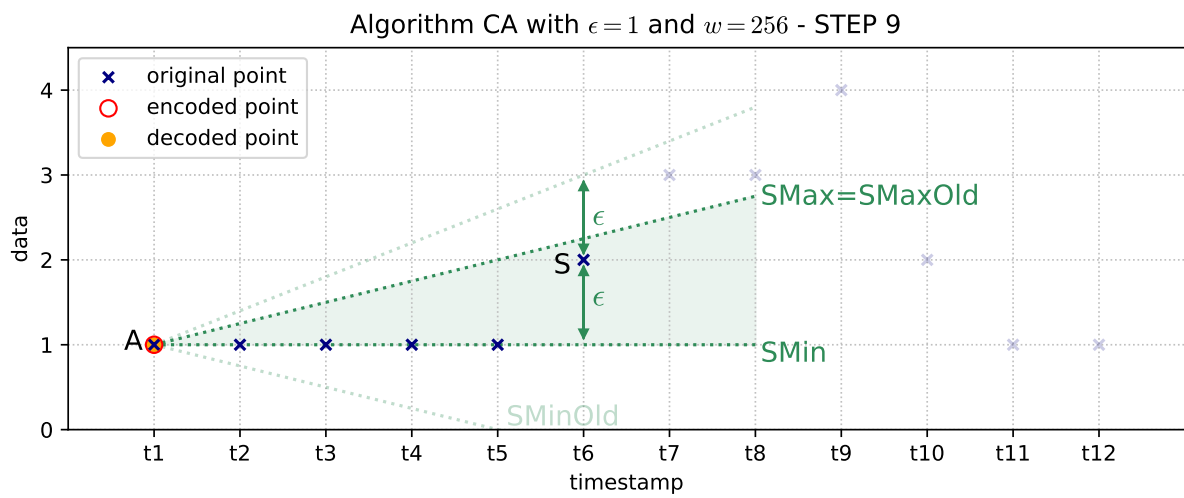


FIGURE 3.35

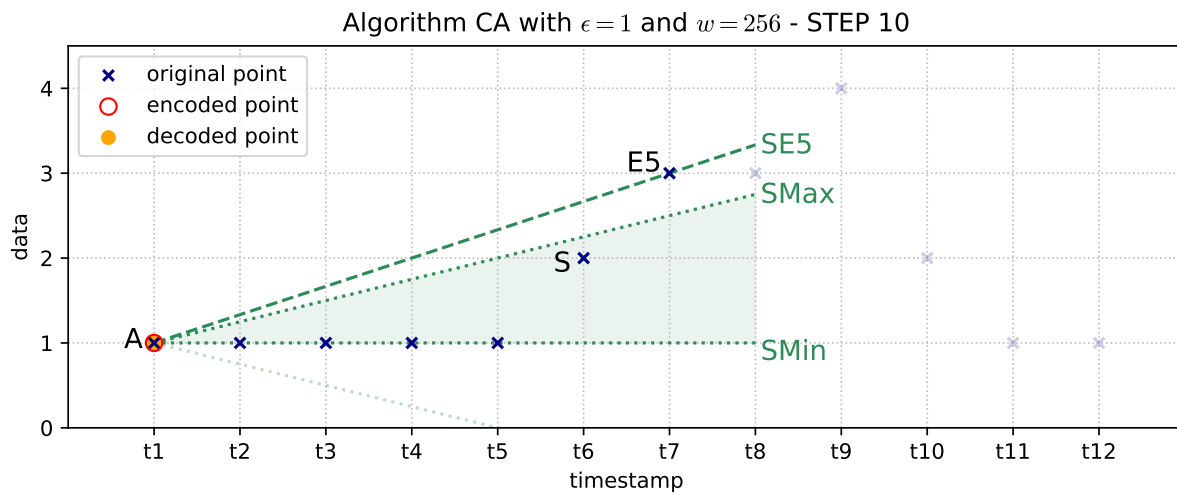


FIGURE 3.36

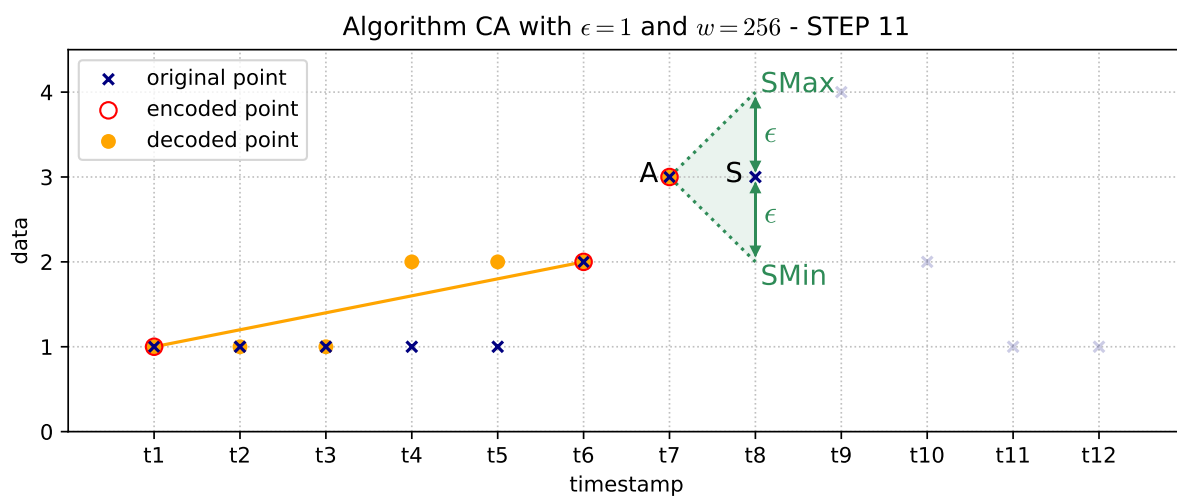


FIGURE 3.37

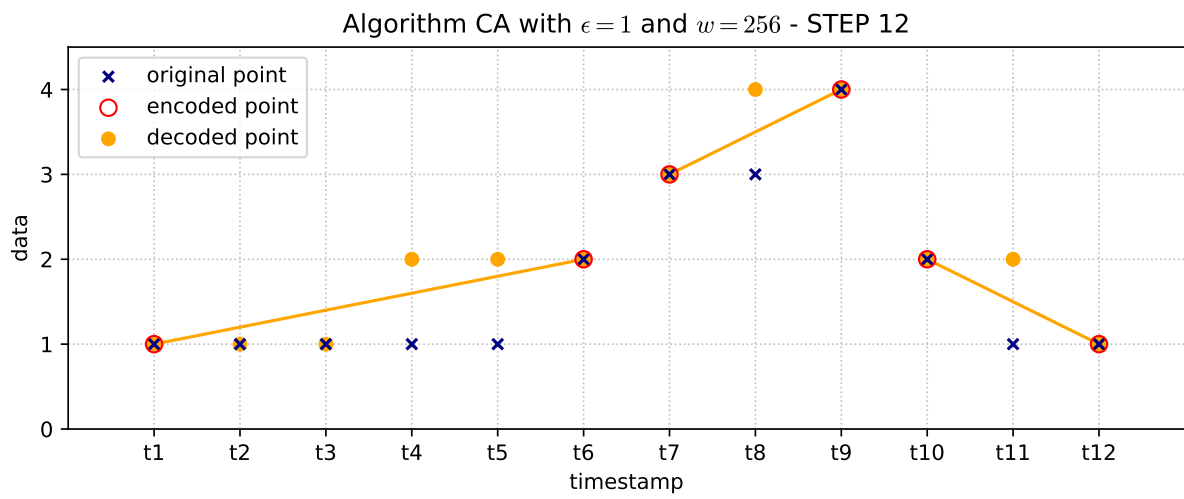


FIGURE 3.38

3.9 Algorithm FR

Algorithm FR [14], 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 3.39 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 encoded
         out: binary file encoded with algorithm FR
          $\epsilon$ : maximum error threshold
          $w$ : maximum window size

1 Create a new window, win
2 foreach entry in column, entry 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  Create a new window, win
14 end

```

FIGURE 3.39: 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 3.40 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 3.40: FR.decode_column_M pseudocode.

In Figure 3.41 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 3.6, 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 (ϵ) 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 3.42.

```

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 3.41: push_points_indexes pseudocode.

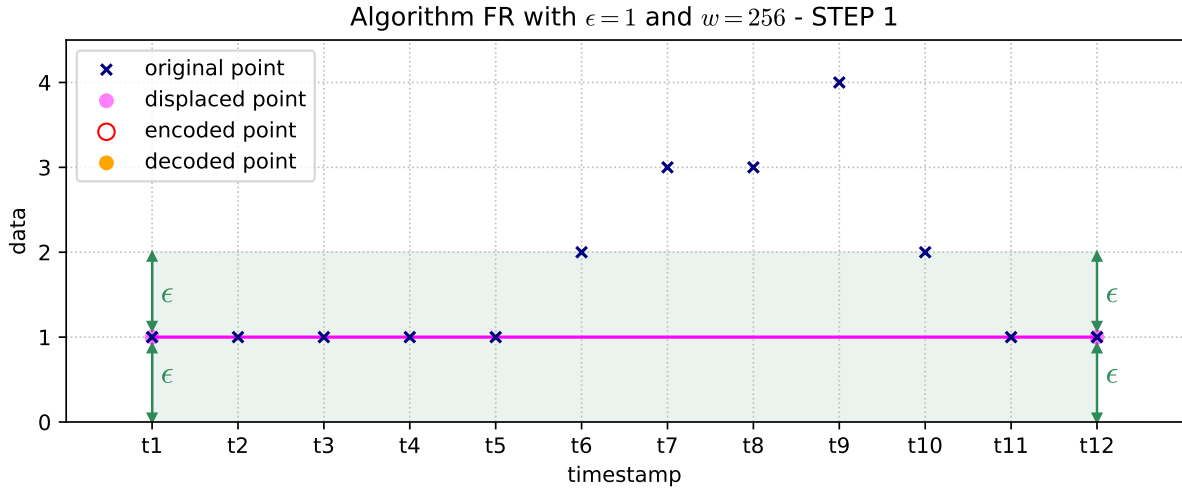


FIGURE 3.42

The condition in line 7 is true, since $0 + 1 < 11$. The condition in line 8 is also true, since as Figure 3.42 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 3.43. 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.

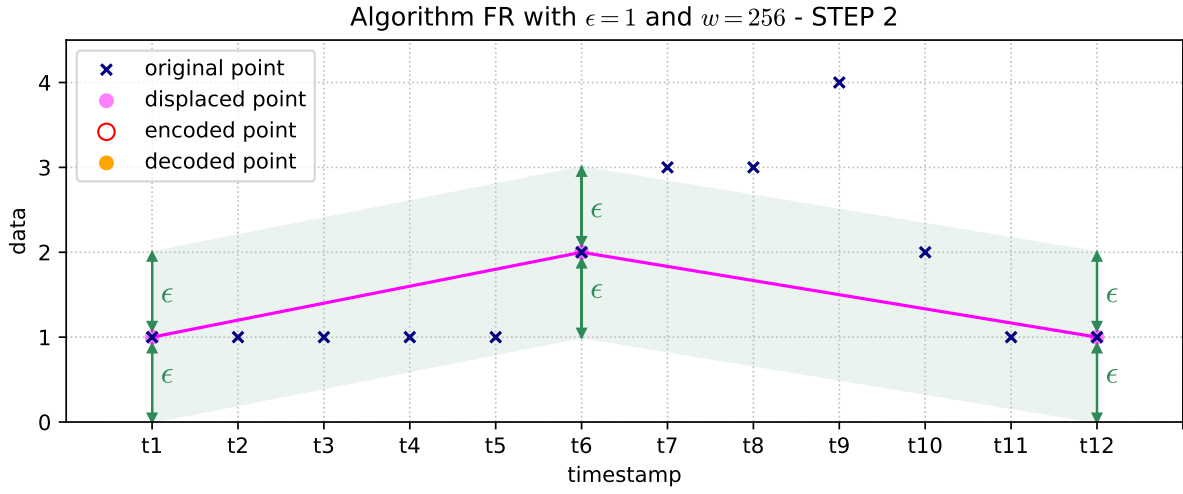


FIGURE 3.43

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 3.43 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 3.44. Observe that all of the points satisfy the error threshold constraint.

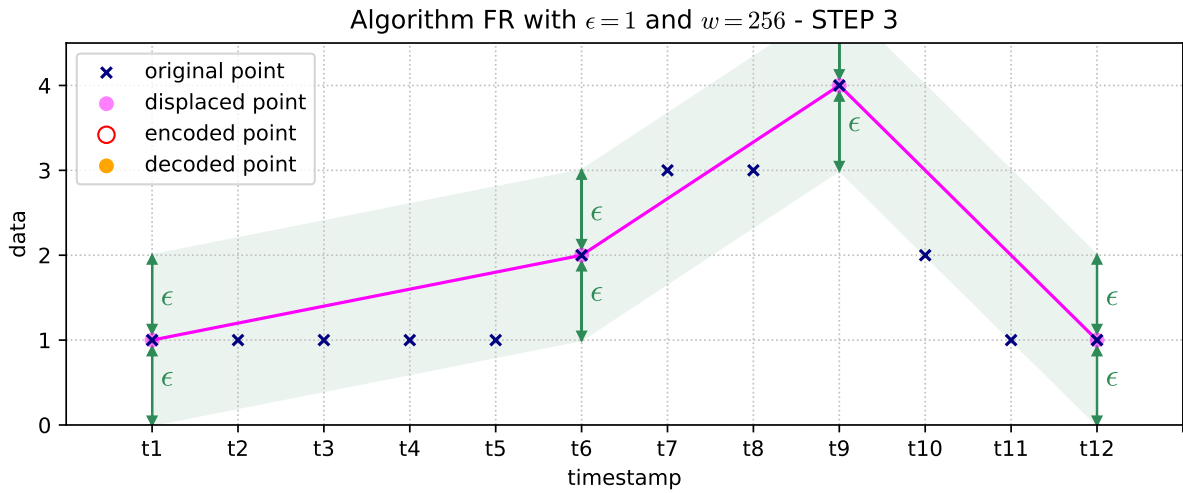


FIGURE 3.44

Finally, in lines 9-12 of the `code_column_M` subroutine, the four endpoints associated to the indexes in $points_indexes$ are encoded. The graph in Figure 3.45 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 3.40, 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.

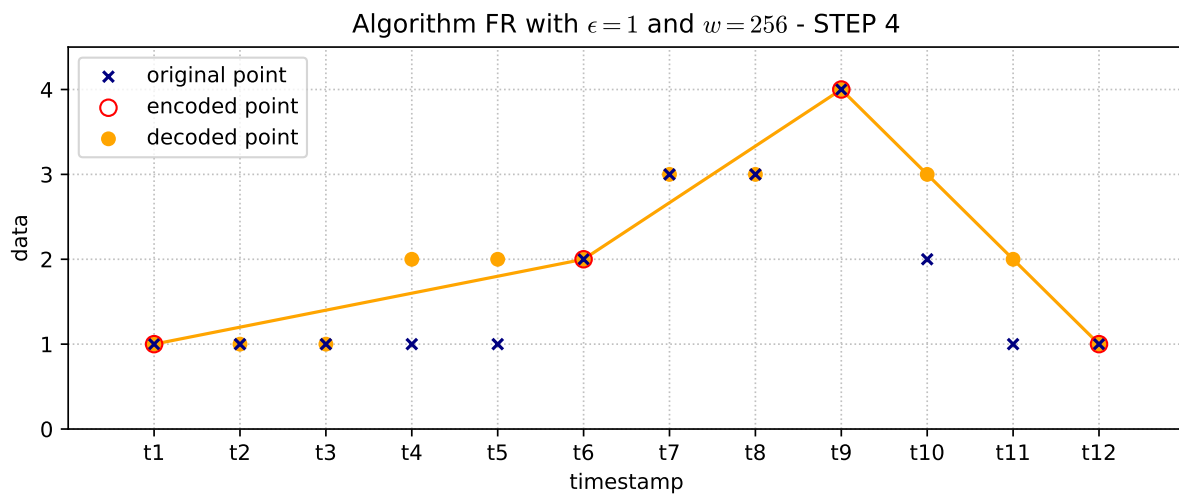


FIGURE 3.45

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

3.10 Algorithm GAMPS

Ver los siguientes documentos:

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

3.11 Other

EL SIGUIENTE PÁRRAFO ESTABA ORIGINALMENTE EN LA SECTION 3.2

All the algorithms are implemented in C++. Our implementations of algorithms PWLH [11], SF [12] and GAMPS [15] reuse part of the source code from the framework cited in [7]¹. The implementations of the remaining algorithms [9, 10, 13, 14] are entirely ours.

EL SIGUIENTE PÁRRAFO ESTABA ORIGINALMENTE EN LA SECTION 3.2.1

The version of the AC algorithm we used in our project is the CACM87 implementation [21, 22]. 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 (recall lines 7-9 in Figure 3.1). 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.

¹The framework is available for download in the following website: <http://lsirwww.epfl.ch/benchmark/>

Chapter 4

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 3, by encoding the various datasets introduced in Chapter 2.

In Section 4.1 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 4.2 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 4.3 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 4.4 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.

4.1 Experimental Setting

We denote by A the set of all the coding algorithms presented in Chapter 3. 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 3.1). 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 2. 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 4.1.1. 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 3.3).

Definition 4.1.2. 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 4.1.3. 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 4.1.4. 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.,

$$CR(c, z, f) = \frac{|c(z, f)|}{|\text{Base}(z, f)|}. \quad (4.1)$$

Notice that smaller values of CR correspond to better performance. Our main goals are to analyze which CAIs yield the smallest values in (4.1) 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 4.1.5. 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

$$RD(c_1, c_2, z, f) = 100 \times \frac{|c_2(z, f)| - |c_1(z, f)|}{|c_2(z, f)|}. \quad (4.2)$$

Notice that c_1 has a better performance than c_2 if (4.2) 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 4.1.3–4.1.5 to datasets, as follows.

Definition 4.1.6. 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 4.1.7. Let z be a data type of a certain dataset d . Let $c \in CI$ be a CAI. We define $|c(z, d)|$ as

$$|c(z, d)| = \sum_{f \in F(d, z)} |c(z, f)|. \quad (4.3)$$

Definition 4.1.8. The *compression ratio* (CR) of a CAI $c \in CI$ for the data type z of a certain dataset d is given by

$$CR(c, z, d) = \frac{|c(z, d)|}{|\text{Base}(z, d)|}. \quad (4.4)$$

Definition 4.1.9. 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

$$RD(c_1, c_2, z, d) = 100 \times \frac{|c_2(z, d)| - |c_1(z, d)|}{|c_2(z, d)|}. \quad (4.5)$$

4.2 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:

- PCA_M against PCA_{NM}
- APCA_M against APCA_{NM}
- CA_M against CA_{NM}
- PWLH_M against PWLH_{NM}
- PWLHInt_M against PWLHInt_{NM}
- GAMPSLimit_M against 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 4.2.1. 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

$$\text{OWS}(a_v, e, z, d) = \arg \min_{w \in W} \left\{ CR(c_{<a_v, w, e>}, z, d) \right\}, \quad (4.6)$$

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 (4.5), 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 4.1 and 4.2 we show the CR and the RD, as a function of the threshold parameter, obtained for two data types of two different datasets. Figure 4.1 shows the results for the data type “SST” of the dataset SST, and Figure 4.2 shows the results for the data type “Longitude” of the dataset Tornado. In Figure 4.1 we observe a large RD favoring the masking variant for all tested algorithms. On the other hand, in Figure 4.2 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.

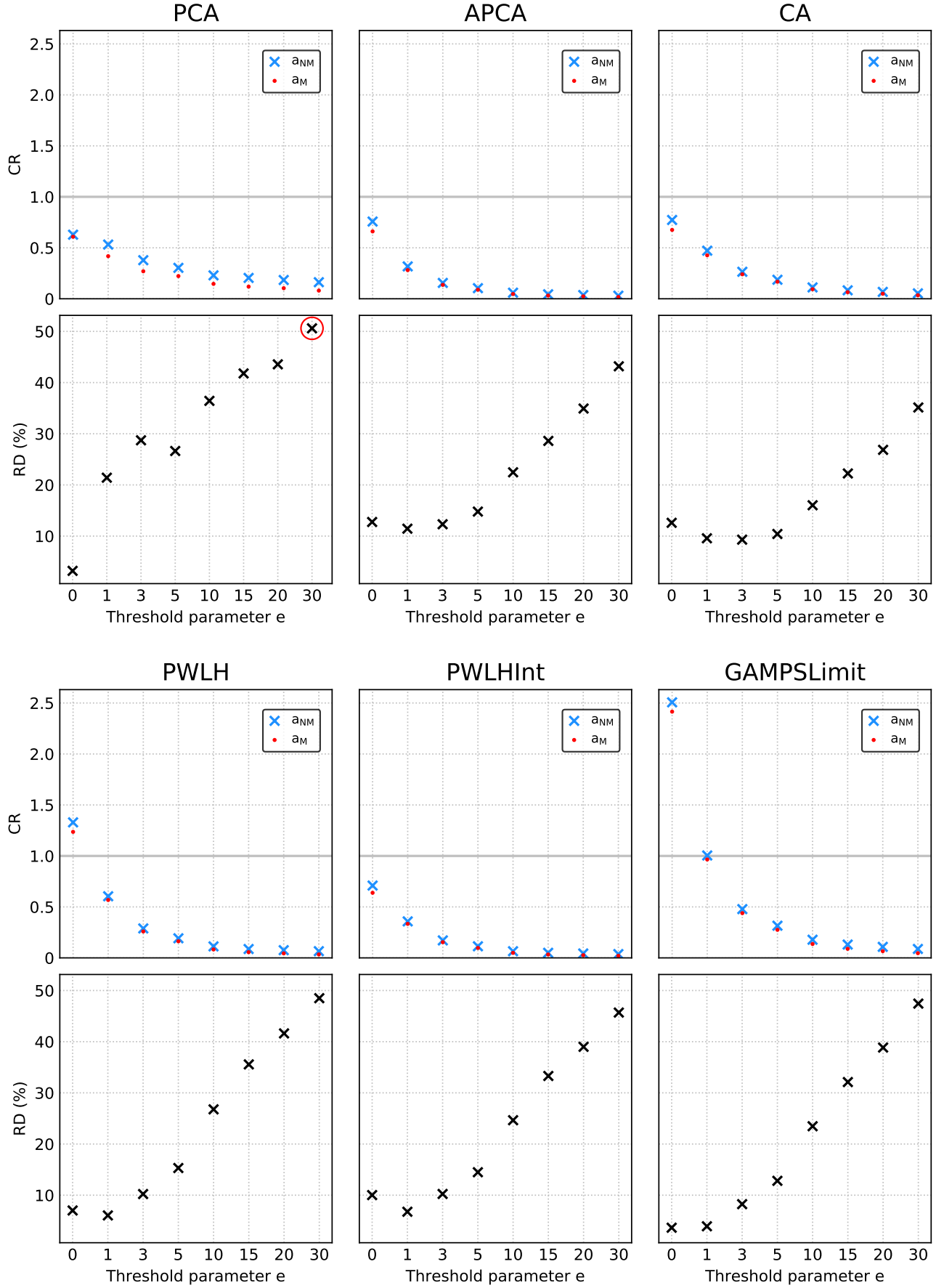


FIGURE 4.1: CR and RD plots for every pair of algorithm variants $a_M, a_{NM} \in A_M$, for the data type “SST” of the dataset SST. In the RD plot for algorithm PCA we highlight with a red circle the marker for the maximum value (50.60%) obtained for all the tested CAIs.

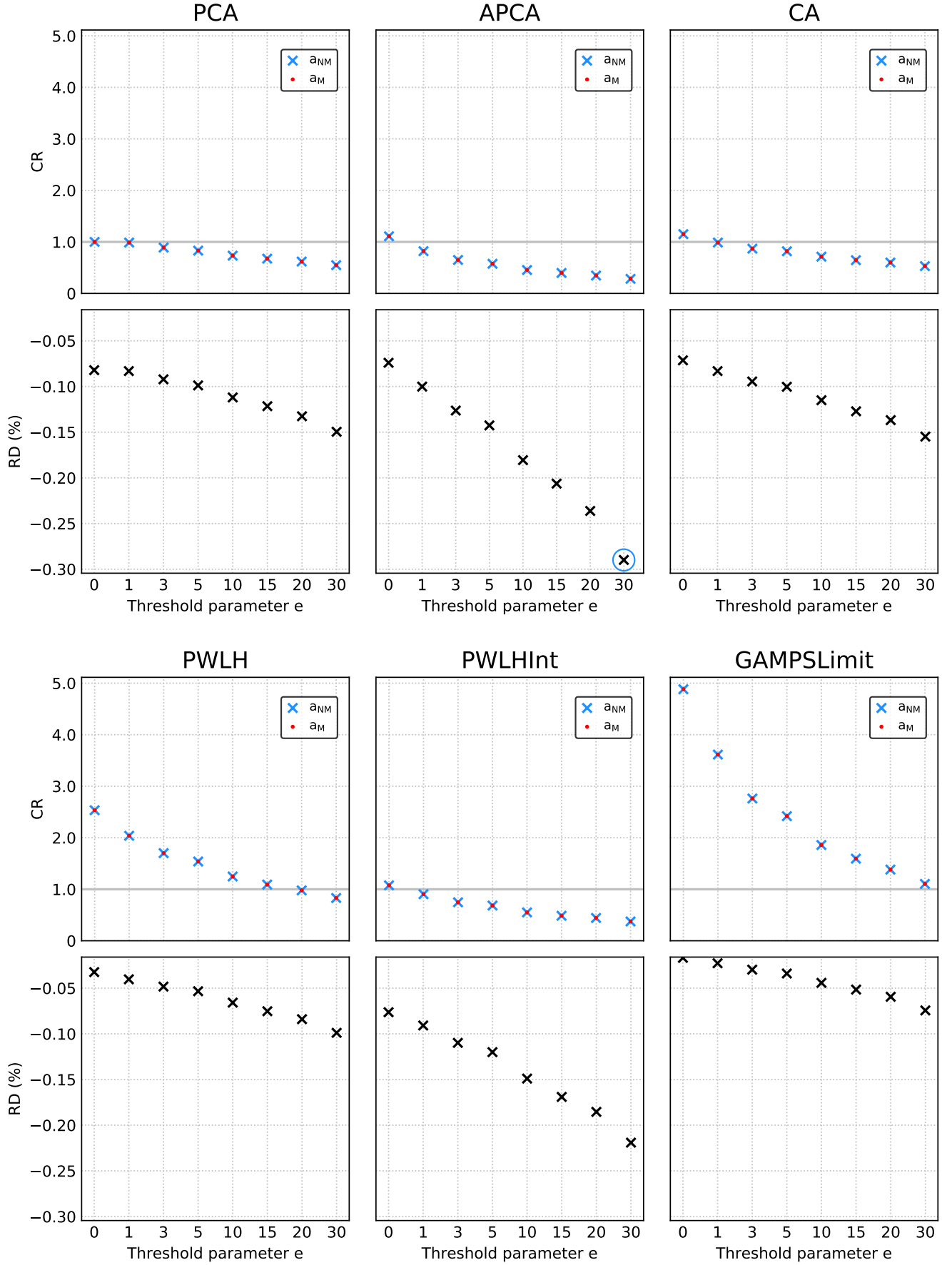


FIGURE 4.2: CR and RD plots for every pair of algorithm variants $a_M, a_{NM} \in \mathcal{A}_M$, for the data type “Longitude” of the dataset Tornado. In the RD plot for algorithm APCA we highlight with a blue circle the marker for the minimum value (-0.29%) obtained for all the tested CAIs.

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 4.1 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 4.1: 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 2.2). 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 4.1, 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 4.2 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 4.1 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 4.2.2. 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.

4.3 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 3.1). Also, we only consider the four datasets that consist of multiple files, i.e. IRKIS, SST, ADCP and Solar (recall this information from Table 2.2). For each file, we compare the compression performance when using the OWS for the dataset, as defined in (4.6), and the LOWS for the file, defined next.

Definition 4.3.1. 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

$$LOWS(a_v, e, z, f) = \arg \min_{w \in W} \left\{ CR(c_{\langle a_v, w, e \rangle}, z, f) \right\}, \quad (4.7)$$

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 (4.2), 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 4.3 and 4.4 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 4.3 shows the results for the file $f = \text{"vwc_1202.dat.csv"}$, and Figure 4.4 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 4.3 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 4.4 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.

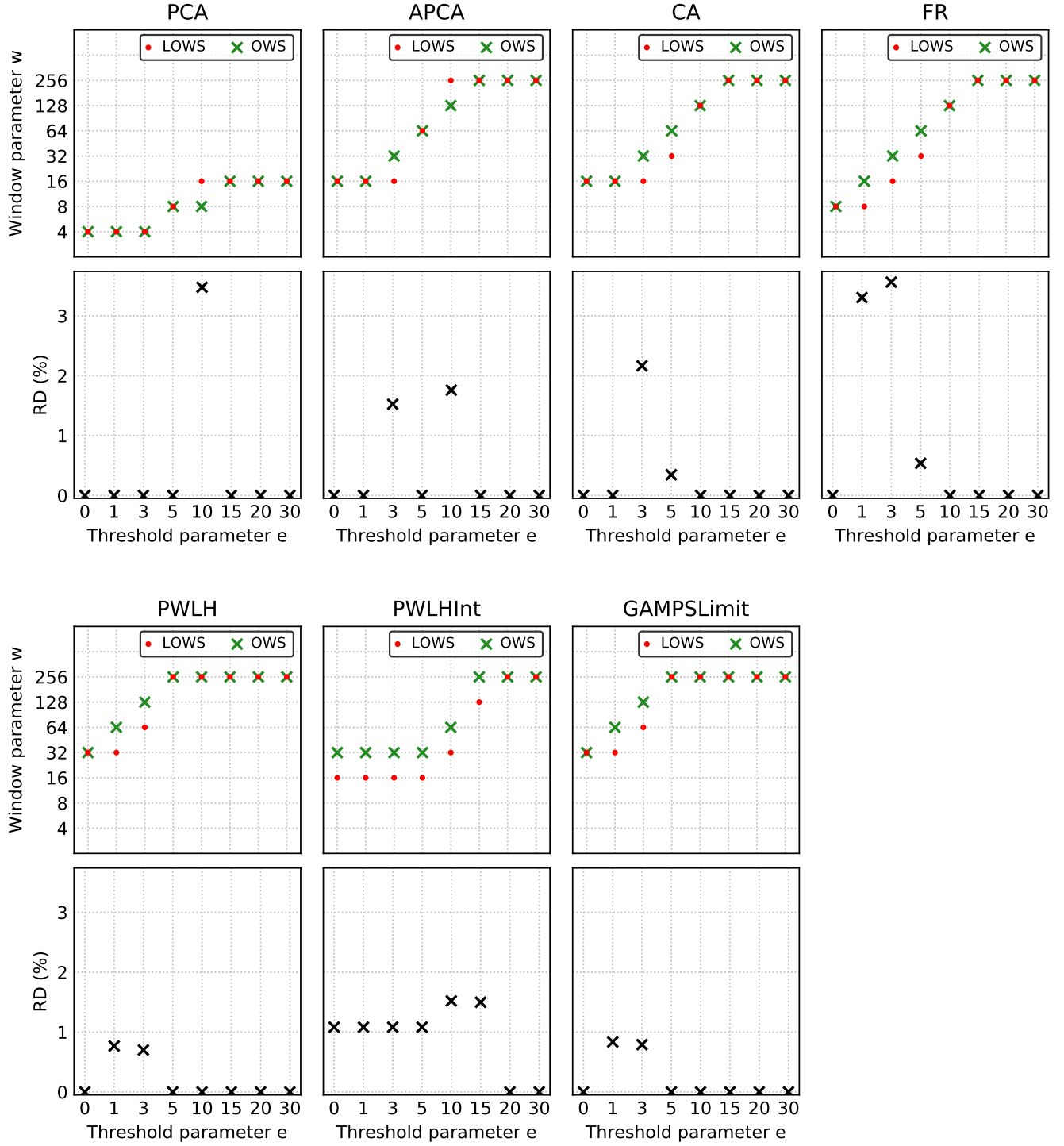


FIGURE 4.3: Plots of w_{global}^* , w_{local}^* , and the RD between $c_{<a_v, w_{global}^*, e>}$ and $c_{<a_v, w_{local}^*, e>}$, as a function of the threshold parameter e , obtained for the data type “VWC” of the file “vwc_1202.dat.csv” of the dataset IRKIS.

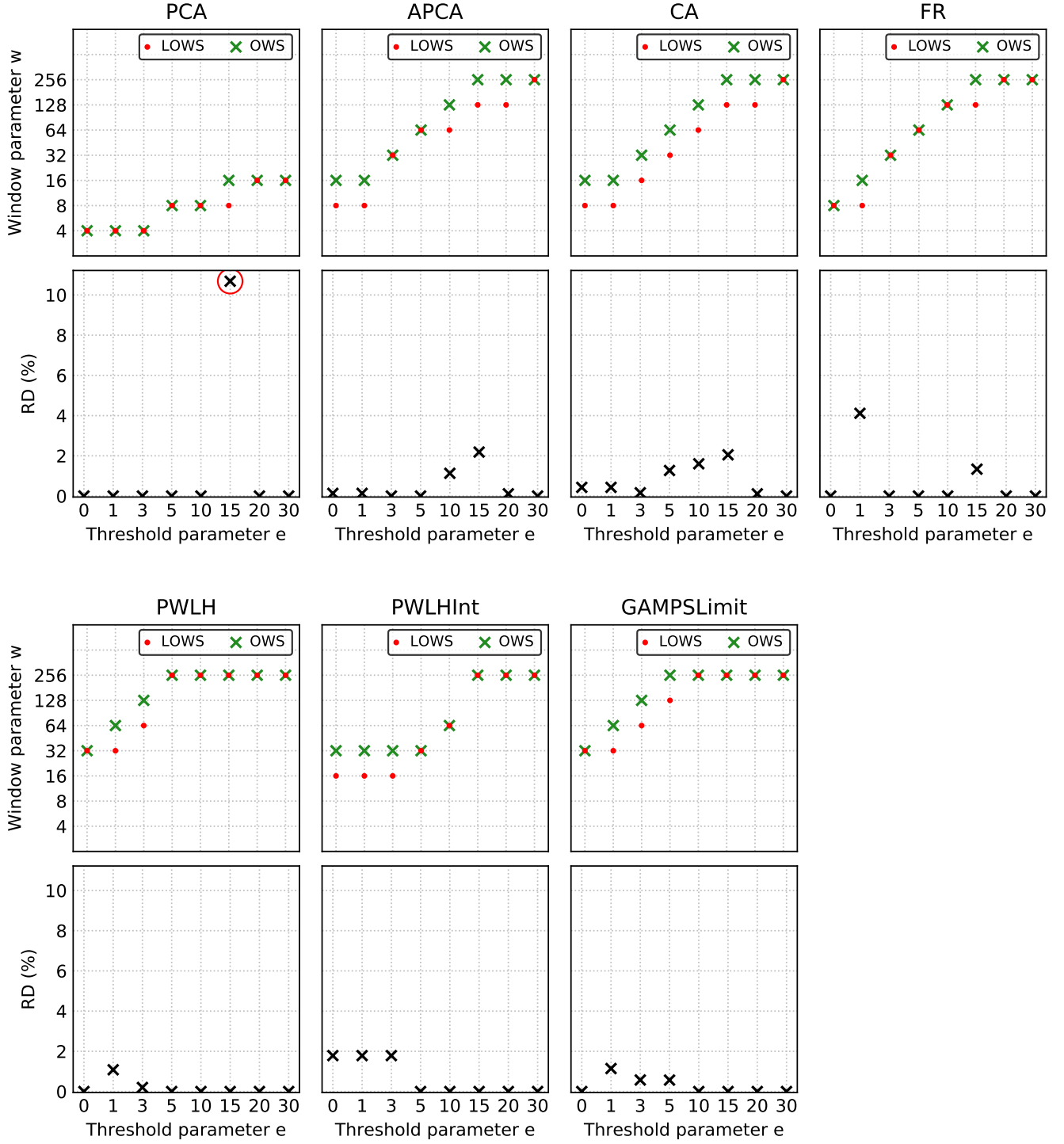


FIGURE 4.4: Plots of w_{global}^* , w_{local}^* , and the RD between $c_{<a_v, w_{global}^*, e>}$ and $c_{<a_v, w_{local}^*, e>}$, as a function of the threshold parameter e , obtained for the data type “VWC” of the file “vwc_1203.dat.csv” of the dataset IRKIS. In the RD plot for algorithm PCA we highlight 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 4.2 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 4.2: 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 4.2 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 4.4 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.

4.4 Algorithms Performance

In this section, we compare the compression performance of the coding algorithms presented in Chapter 3, by encoding the various datasets introduced in Chapter 2. 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 4.1.6–4.1.9. Following the results obtained in sections 4.2 and 4.3, 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 4.2.1).

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 (4.4), 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 4.4.1. 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.,

$$c^b(z, d, e) = \arg \min_{c \in CI} \left\{ CR(c_{\langle a_v, w_{global}^*, e \rangle}, z, d) \right\}. \quad (4.8)$$

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 4.5 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.

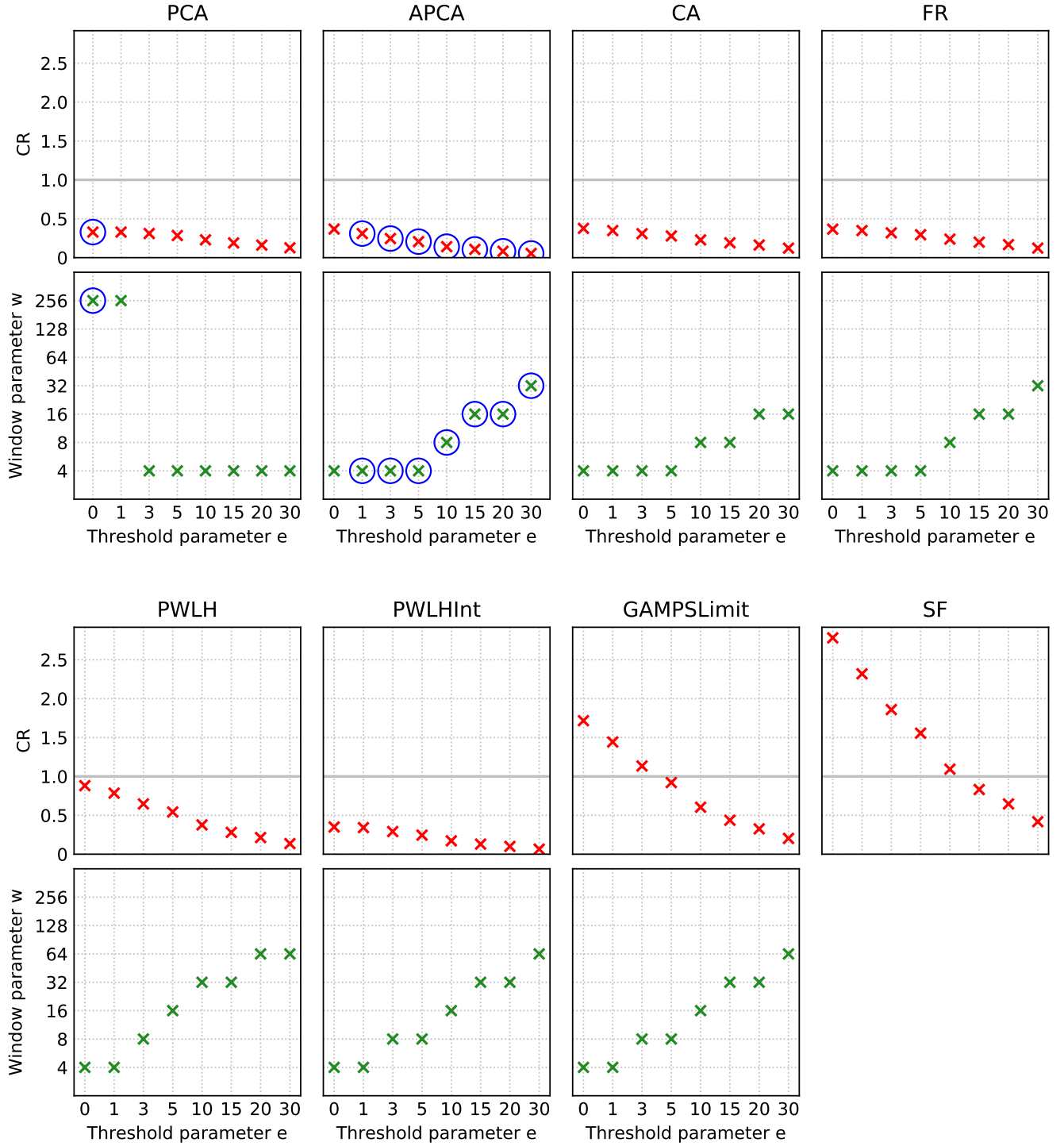


FIGURE 4.5: CR and window size parameter plots for every algorithm, for the data type “SST” of the dataset ElNino. For each threshold parameter $e \in E$, we use blue circles to highlight the markers for the minimum CR value and the best window size parameter (in the respective plots corresponding to the best algorithm)

Table 4.3 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 4.5. 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 4.3: 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 4.4.2. The *maximum RD* (maxRD) of a coding algorithm $a \in A$ for certain threshold parameter $e \in E$ is given by

$$\text{maxRD}(a, e) = \max_{z, d} \left\{ \text{RD}(c^b(z, d, e), c_{<a_v, w_{global}^*, e>}) \right\}, \quad (4.9)$$

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 4.3, 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 4.4.3. The *minmax RD* (minmaxRD) for certain threshold parameter $e \in E$ is given by

$$\text{minmaxRD}(e) = \min_{a \in A} \left\{ \text{maxRD}(a, e) \right\}, \quad (4.10)$$

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 4.4 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 4.4: $\text{maxRD}(a, e)$ 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 4.3 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.

4.4.1 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 4.5 summarizes the compression performance results obtained by gzip and the other evaluated coding algorithms, for each data type of each dataset. Similarly to Table 4.3, 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 4.5: 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 4.4, Table 4.6 shows the $\text{maxRD}(a, e)$ 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 correspondent to the $\text{minmaxRD}(e)$ value is highlighted.

We observe that, for every e , the minmaxRD 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 minmax coding algorithm when $e \in [0, 1]$, and in both cases the minmaxRD values are rather high, i.e. 26.58% and 47.98%, respectively. APCA remains the minmax coding algorithm for every $e \geq 3$, but its minmaxRD 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 4.6: $\text{maxRD}(a, e)$ 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.

4.5 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 4.6.

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.

4.6 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] N. Wever. IRKIS Soil moisture measurements Davos. SLF. <https://doi.org/10.16904/17>, 2017. [Accessed 9 October 2020].
- [2] N. Wever, F. Comola, M. Bavay, and M. Lehning. Simulating the influence of snow surface processes on soil moisture dynamics and streamflow generation in an alpine catchment. *Hydrol. Earth Syst. Sci.*, 21:4053–4071, <https://doi.org/10.5194/hess-21-4053-2017>, 2017.
- [3] NOAA - TAO Data Download. https://tao.ndbc.noaa.gov/tao/data_download/search_map.shtml, 2016. [Accessed 9 October 2020].
- [4] SolarAnywhere - Data. <https://data.solaranywhere.com/Data>, 2020. [Accessed 9 October 2020].
- [5] El Nino Data Set. <https://archive.ics.uci.edu/ml/datasets/El+Nino>, 1999. [Accessed 9 October 2020].
- [6] Kaggle - Storm Prediction Center. <https://www.kaggle.com/jtennis/spctornado>, 2016. [Accessed 9 October 2020].
- [7] 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.
- [8] 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.
- [9] I. Lazaridis and S. Mehrotra. Capturing Sensor-Generated Time Series with Quality Guarantees. *Proc. ICDE*, pages 429–440, 2003.
- [10] 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.
- [11] C. Buragohain, N. Shrivastava, and S. Suri. Space Efficient Streaming Algorithms for the Maximum Error Histogram. *Proc. ICDE*, pages 1026–1035, 2007.
- [12] 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.
- [13] G.E. Williams. Critical Aperture Convergence Filtering and Systems and Methods Thereof. *US Patent 7,076,402*, Jul. 11, 2006.

- [14] J.A.M. Heras and A. Donati. Fractal Resampling: time series archive lossy compression with guarantees. *PV 2013 Conference*, 2013.
- [15] 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.
- [16] J.J. Rissanen. Generalized Kraft Inequality and Arithmetic Coding. *IBM Journal of Research and Development*, 20(3):198–203, 1976.
- [17] Arithmetic Coding + Statistical Modeling = Data Compression. <https://marknelson.us/posts/1991/02/01/arithmetic-coding-statistical-modeling-data-compression.html>, 1991. [Accessed 9 October 2020].
- [18] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley-Interscience, 2nd edition, 2006.
- [19] R.E. Krichevsky and V.K. Trofimov. The performance of universal encoding. *IEEE Transactions on Information Theory*, 27(2):199–207, 1981.
- [20] R. Graham. An Efficient Algorithm for Determining the Convex Hull of a Finite Planar Set. *Information Processing Letters*, 1:132–133, 1972.
- [21] I.H. Witten, R.M. Neal, and J.G. Cleary. Arithmetic Coding for Data Compression. *Communications of the ACM*, 30(6):520–540, 1987.
- [22] Data Compression With Arithmetic Coding. <https://marknelson.us/posts/2014/10/19/data-compression-with-arithmetic-coding.html>, 2014. [Accessed 9 October 2020].
- [23] GNU Gzip. <https://www.gnu.org/software/gzip/>, 2018. [Accessed 9 October 2020].
- [24] GNU Gzip Manual - Invoking gzip. https://www.gnu.org/software/gzip/manual/html_node/Invoking-gzip.html, 2018. [Accessed 9 October 2020].
- [25] Solomon W. Golomb. Run-length encodings (Corresp.). *IEEE Transactions on Information Theory*, 12(3), 1966.