

Journal of Integrated Information Management

Vol 9, No 2 (2024)

Jul-Dec 2024



**Information: A physical reality or a humanly tool?
From the model order to the appropriate number of
clusters**

Markos Dendrinou

Copyright © 2024, Markos Dendrinou



This work is licensed under a [Creative Commons Attribution-NonCommercial 4.0](https://creativecommons.org/licenses/by-nc/4.0/).

To cite this article:

Dendrinou, M. (2024). Information: A physical reality or a humanly tool? From the model order to the appropriate number of clusters. *Journal of Integrated Information Management*, 9(2), 7-13. Retrieved from <https://ejournals.epublishing.ekt.gr/index.php/jiim/article/view/39839>

Information: A physical reality or a humanly tool? From the model order to the appropriate number of clusters

Markos Dendrinios

Professor of Informatics and History & Philosophy of Science, Department of Archival, Library and Information Studies, University of West Attica, Athens, Greece

mdendr@uniwa.gr [ORCID: 0000-0001-5675-3069]

Article Info

Article history:

Received 06 May 2024

Received in revised form 12 June 2024

Accepted 25 June 2024

<http://dx.doi.org/10.26265/jiim.v9i2.39839>

Abstract:

This paper is a presentation of two important types of information regarding natural signals and groups of relative things. The measure of the first information type is the order of the system that produces the signal, while the second one is the appropriate number of distinct clusters for the most effective classification of a certain group into narrower and more representative subgroups. Firstly a review is given concerning the two approaches, and then a certain method is proposed for the identification of the correct number of clusters to be used in K-Means clustering process. The two information measures, model order and number of clusters, can be considered as two equivalent views of an inherent natural element, an objective order behind any physical system.

Index Terms — information, information measure, information detection, model order, K-Means, clustering.

I. INTRODUCTION

Heisenberg (1927) has formulated the principle of uncertainty. According to this principle, the combined accuracy of measurement of the position and the momentum (velocity) of a particle has a lower bound, which is of the order of magnitude of the Planck constant h . This principle constitutes an alarm against the traditional belief of the unavoidable scientific progress in the description of the reality. The empiricists emphasize that Heisenberg's principle of uncertainty concerns mainly nanoscale measurements and its power seems to be eliminated at the macroscopic level. This fact does not override the ultimate scientific truth that the humanly observations and measurements have certain limits and we must face the reality in other terms. Heisenberg and the radical group of the quantum physicists unsettled the scientific background even more, as they denied the physical reality at all.

Thus a fundamental question has been set: Do we have the right to assign a value to the momentum of an electron just before its final measurement? Is it real? According to Heisenberg it is not. Before the final measurement, the best we can attribute to the electron is some unsharp or fuzzy momentum. The ontological sense of the elementary units of matter has been overturn¹.

If we add to the whole image the role of the observer in the double slit experiment, where the nature and the effects of the beam of electrons shot against the slit changes depending on the presence or not of an observer, then, indeed, we can say that reality has lost all its supposed compactness becoming a “smoke of probabilities”. We must also take into account the complex not exclusive nature of the matter, which is particle-like and wave-like simultaneously, as well as the quantum entanglement between any couple of particles, where the measurement of any characteristic of the first affects in no time the value of the second. These properties were the stimulus for quantum computing where the basic unit of information, the qubit has not a distinct certain value 0 or 1 but it is a superposition of both of them, obtaining a value only in the case of its use, affecting also any other qubit entangled with it.

From all the above we can see that the structure of the universe as well as the nature of the information has radically changed. I'll try in the subsequent sections to find a suitable place for an ontological structure of information within this extraordinary contemporary view of the reality. In such a perspective the information could claim the role of the ultimate element of reality substituting the particle-like or wave-like matter.

Information has a lot of significances: a measure of expectancy towards the most economical representation (information entropy), a language system for human communication along with intentionally created symbolical languages (propositional and predicate calculus), a number of harmonics (resonances) of a natural signal related to the order of the underlying system, classification of an object to a predefined or unknown number of clusters. The last two cases will be the subject of this paper, in an attempt to find any common ground between them and a way to shed a light to the open question of the objectivity or human

¹ <https://plato.stanford.edu/entries/qt-uncertainty/>

dependent nature of the information. Model order detection and clustering processes will be used as the tools to this end.

II. INFORMATION VS NOISE

In this section we'll make a review in the literature concerning order detection of synthetic harmonic signals as well as natural signals, such as speech signals and respiratory signals. The order of these signals is meant as their harmonic content. In the case of the natural signals this is an objective information indicator, corresponding to the anatomical characteristics of the system that produces them (phonetic tube, lung structure etc.). In the case of synthetic harmonic signals the order is twice their number of harmonics; the order detection is problematic mainly in the case of a signal embedded in white or colored noise and it proves to be of low success when the Signal to Noise Ratio (SNR) is low. The general idea is to reduce the dimensionality of a signal retaining only the dimensions that contain the real information part through applying Principal Component Analysis [3].

In this frame SVD methodology can be applied in the data or autocorrelation matrix of the signal in order to find the singular values (a generalization of eigen-values) of the matrix and remove all the unnecessary values which are supposed to correspond to the noise. The required order is the number of the remaining singular values, which consist the information subspace. The identification of this number is crucial, since this number shows the orthogonal axes into which the real information is decomposed, that is the signal subspace against the noise subspace. An analysis by synthesis procedure based on the SVD methodology has been proposed in [1] for synthetic signals and in [8] for speech signals, which has been a basic reference for many articles about order detection, speech enhancement and noise cancellation since then.

Bakamidis, Dendrinis and Carayannis [1], in the frame of decomposing the synthetic harmonic signals through SVD Principal Component Analysis, defined a new criterion NEE (standing for Noise Error Estimation), superior to the Rao criterion ([12]), aiming at detecting the number of harmonics in the presence of noise [1]. Consecutive reconstructions are performed and the resulting error power is compared to the noise variance in order to get the best approximation of the original non corrupted signal. The number of the singular values corresponding to a reconstruction error power as close as possible to the noise variance gives the parsimonious order. The existence of such a criterion is important for both high quality reconstruction and accurate spectral analysis. Various spectral estimation techniques used on the reconstructed signal give the possibility to retrieve harmonics in highly noisy environment along with availability of very short data lengths.

In this frame, a number of successively reconstructed signals, derived through this analysis by synthesis procedure, were used for determination of the harmonics order of synthetic signals in hard conditions: limited number of samples, low SNR, and close frequencies. The

introduced criterion shows not only the number of the harmonics (complex sinusoids) but also the best order for reconstruction (leading to a signal as close to the pure one as possible). This methodology has been applied in the field of enhancement of noisy speech with great success.

Some indicators of the success rates of the proposed method are the following: Case A. Signal length 100, SNR 0 db, sinusoidal frequencies: $0.4 f_0$, $0.43 f_0$ (f_0 is the Nyquist rate): Success rate: 96%. Case B. Signal length 100, SNR -3 db, sinusoidal frequencies: $0.4 f_0$, $0.43 f_0$: Success rate: 80%. This proved to be a significant improvement compared to the 83% and 53% of Rao criterion (Rao et. al., 1988), in spite of the shorter used data records.

NEE criterion has been also applied in real conditions, as in speech enhancement from noise and detection of Obstructive Sleep Apnea (OSA) in medicine. In both cases the detection of the exact order of the signal is the key factor for the correct reconstruction of the signal and the abolishment of noise along with the minimum distortion of the initial signal.

Dendrinis, Bakamidis and Carayannis [8] proposed a speech enhancement technique, based on principal component analysis and a new criterion for the selection of the parsimonious number of real signal components aiming at noise-free signal regeneration [8]. Both isolated phoneme and continuous speech experiments were presented. The results have been evaluated by informal listening and SNR computations, which show that the methodology had an improved performance compared to other widely used methods. The estimated number of singular values to be retained is supposed to keep all the necessary information content of the speech signal, while the omitted ones correspond to the additive noise.

In this context various phoneme and phrases embedded in noise have been reconstructed after recovering the correct model order of the natural signal.

As you can see in Fig. 1 (a), the order of a noisy pronounced phoneme [e] was selected where the NEE measure takes its lower value (order=9). The reconstructed signal corresponding to this value is the signal with the greater segmental SNR, showing the best possible noise elimination. Similarly, in (b) the order of a noisy pronounced phoneme [n] was selected where the NEE measure takes its lower value (order=6).

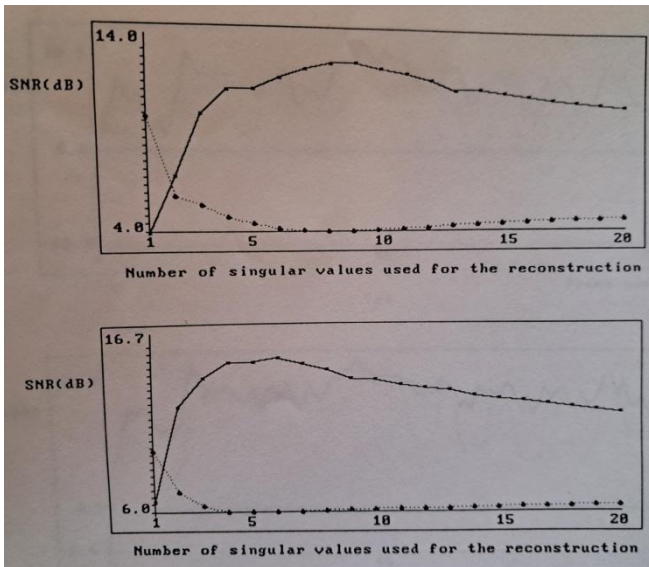


Fig.1. (a) The segmental SNR's (continuous line) and the NEE criterion values (dotted line) for the phoneme [e] embedded in noise 10 dB, (b) The same for the phoneme [n].

Next figure shows the recovered spectral characteristics of a Greek dictated phrase, which had been hidden due to the noise. The phrase was divided in many overlapping frames, each of which was analyzed. After selecting the parsimonious order for each frame, a reconstructed signal is derived. The correct order selection for each frame is of great importance in this process.

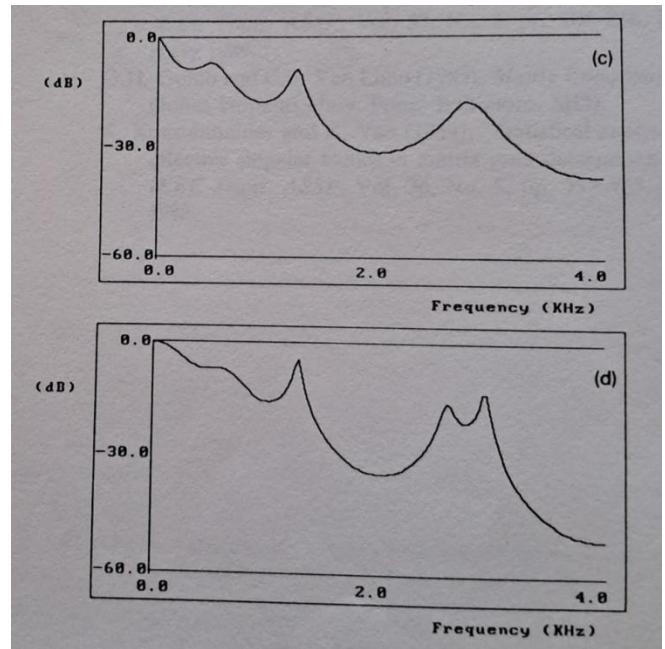
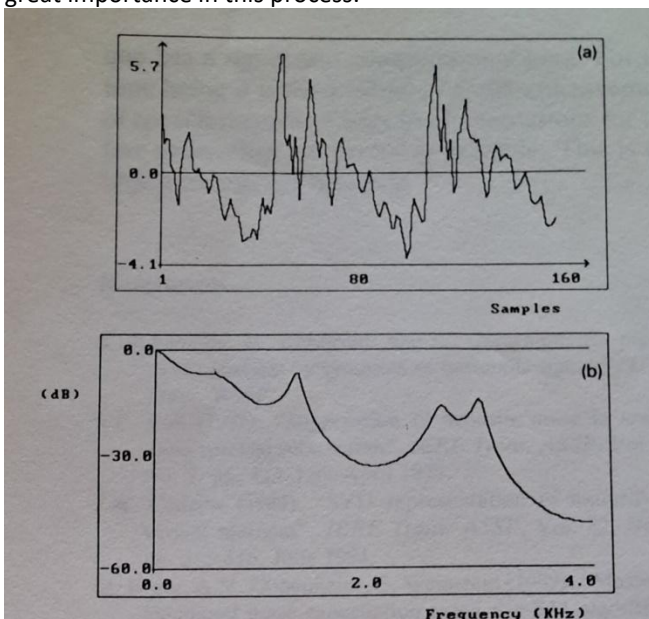


Fig. 2. (a) An examined signal frame [No. 21] of a dictated sentence in Greek ('i mana mou emene me mena') embedded in noise with global SNR = 20 dB, (b,c,d) The LPC (Linear Predictive Coding) spectral densities of the original, the noisy and the reconstructed signals, respectively.

Since then, various alternative approaches, taking into account more or less the NEE technique, have been suggested, which try to enhance further the order detection efficacy. Besides, a lot of speech enhancement methods based more or less in the above methodology have also been suggested in the meantime.

Wang et al (2015) [15] propose an efficient framework and algorithm for the one dimensional harmonic retrieval problem in the case of either additive colored Gaussian or non-Gaussian noise, in the hard condition of very closely spaced harmonic signals in the frequency domain. The wavelet packet (WP) method to the blind source separation (BSS) based harmonic retrieval model is utilized.

Tong et al. (2017) [14] study the problem of parallel waveform enhancement via the multi-sensor fusion technology. They represent the observed multiple noisy observations as a 3-D tensor, and propose two novel approaches in the time domain, i.e. the transforming and filtering (TAF) approach and the direct multidimensional filtering (DMF) approach, for parallel waveform recovery and interference suppression. The system can produce an estimate of the clean waveform in each sensor channel simultaneously, as it is implied by the term "parallel". In the TAF approach the observed tensor is transformed into a different domain where the noise can then be filtered by discarding the insignificant coefficients. Noise reduction is achieved in the DMF approach directly by applying multidimensional filtering on the observed tensor. Both approaches are "blind" since they do not require a priori frequency responses between the desired source and distributed sensors.

Deppisch et al. (2023) [9] propose a similar subspace method that decomposes SRIRs (Spatial room impulse

responses) into a direct part, which comprises the direct sound and the salient reflections, and a residual, to facilitate enhanced analysis and rendering methods by providing individual access to these components. The proposed method is based on the generalized singular value decomposition and interprets the residual as noise that is to be separated from the other components of the reverberation. Large generalized singular values are attributed to the direct part, which is then obtained as a low-rank approximation of the SRIRs.

Balasubramanian et al. (2023) [2] suggest a combined method for speech enhancement, based on both speech cochleagram and visual cues using Audio-Visual Multichannel Convolutional Neural Network (AVMCNN). Several researchers have already shown that speech enhancement using visual data as an additional input along with the audio data is more effective in minimizing the acoustic noise present in the speech signal. This work proposes a novel CNN-based audio-visual Ideal Ratio Mask (IRM) estimation model.

Christensen et. al. (2016) [6] investigate the potential performance of generalized subspace filters for speech enhancement in cocktail party situations with very poor signal/noise ratio, e.g. down to -15 dB. Performance metrics such as output signal/noise ratio, signal/distortion ratio, speech quality rating and speech intelligibility rating are mapped as functions of two algorithm parameters. The paper results to a recommended trade-off between noise, distortion and subjective performances.

Dendrinios (1998) [6] has used spectral characteristics of the respiratory signals in order to discriminate between normal (just breathing or snoring) and pathological (apnea) cases. Research has been held for the selection of the suitable features of the examined signals. The features which were proved to be information distinctive were the spectral energy between 0 and 700 Hz and the one between 700 and 1400 Hz. The spectral ration between them takes high values in the cases of apneas and pauses, whereas it takes much lower values during breathing or snoring periods.

Steinhaouer et al. (2019) [13] present a system to detect symptoms of allergic rhinitis remotely by using uttered speech and by exploiting its specific spectral characteristics. Based on the principles of adaptive modelling and fundamental frequency variations (jitter) as well as speech analysis by means of acoustic models, the proposed technique achieves an efficient classification of patients from uttered speech. A Singular Value Decomposition based iterative approach is used for the accurate estimation of the jitter and Hidden Markov Models are implemented to model the 32 phonemes. The final decision is derived by optimally combining the individual estimates providing a tool for the automatic diagnosis of allergic rhinitis.

III. INFORMATION AS A CLASSIFICATION TOOL

The model order, conceived as the order of a system that produces natural signals, such as speech or respiratory signals, is an indicator of an objective information,

independent of the humanly intervention. This can be considered as information hidden in nature, which, under the appropriate approaches, as those developed in the previous section, can be revealed.

In this section we'll deal with a continuous implicit or explicit process in the frame of human understanding of our environment: the classification of things or events in categories aiming at an easier confrontation of the 'information' that incessantly shoots us. Data classification gives the human the first phase of the sense of the things, as they are classified into categories (or classes) according to some carefully selected features. The origin of this endeavor goes back to ancient times, in the philosophical context of platonic ideas, which were seen as representatives of groups of things with similar characteristics. There is a great discussion about the objectivity or subjectivity of the process of categorization, which can be thought of as a part of the open question of the nature of the information itself.

Data analysis and information generation is a critical component of intelligence. The formulation of questions, collection and preparation of data, analysis and interpretation which in turn may lead to responses by the intelligent agent are not restricted to humans. Several living organisms have the ability to obtain their environment, extract conclusions, and produce responses. They may also predict events or count or cluster similar objects.

In the frame of classification there is a great branch of statistical informatics offering us tools for supervised or unsupervised clustering, that is division of a number of given vectors of data to a number of clusters. The data here are the features of the objects to be discriminated and they are either numerical or nominal. The number of clusters (classes) where the objects will be classified is critical and it reminds the model order selection problem, discussed in the previous section. One widely used method is the trade-off between minimizing the within-cluster sum of squares (a measure of how tight each cluster is) and maximizing the between-cluster sum of squares (a measure of how separated each cluster is from the others)².

In this context we are seeking a method for absolutely automatic unsupervised clustering leading to the estimation of model order (the parsimonious number of clusters). The results are compared to the experts' estimation of that order, since we suppose that an experts group has a more accurate perception of a certain area; for example a botanist may conclude differently to iris flowers attributes, due to a 'positive bias' extracted from his exhausting measurements, whilst a non-expert may not. In this way objective and subjective information detection seem to mutually converge.

The method proposed in this paper is the consecutive application of K-means clustering algorithm, through increasing one by one the number of the clusters. The used environment is the open source data mining platform WEKA, We compute in each case the difference of the 'within cluster sum of squared errors' SSE value from the

²

<https://towardsdatascience.com/10-tips-for-choosing-the-optimal-number-of-clusters-277e93d72d92>

corresponding one in the previous case (current number of clusters-1). Then the rhythm of decrease (difference of differences) is estimated as a discrete measure of second derivative of SSE. The maximum value of the sequence of the 'rhythm of decrease' values gives the optimal number of clusters.

Two classification cases are presented. The one case concerns a group of iris flowers characterized by 4 metrical features of their petals and sepals. The SEE criterion gives the true number of categories-clusters, that is 3, the same with the a-priori classification of the flower instances given by the plant scientists into Iris Setosa - Iris Versicolour - Iris Virginica. The second case concerns a group of image segmentation data extracted from 7 distinct categories of outdoor images: brickface, sky, foliage, cement, window, path, grass. The SEE criterion gives again the true number of categories-clusters, that is 7.

A. Iris Plants Database

- **Creator:** R.A. Fisher.
- **Donor:** Michael Marshall (MARSHALL%PLU@io.arc.nasa.gov), July, 1998.
- **Number of Instances:** 150 (50 in each of three classes).
- **Number of Attributes:** 4 numeric attributes and the class.
- **Attribute Information:** (1) sepal length in cm, (2) sepal width in cm, (3) petal length in cm, (4)] petal width in cm, (5) class: Iris Setosa - Iris Versicolour - Iris Virginica.
- **Missing Attribute Values:** None.

You can see below the data extracted after the consecutive K-means applications.

Table 1. Iris flower. K-Means clustering in WEKA environment. (a) Number of clusters. (b) SSE: within cluster sum of squared errors. (c) Difference between consecutive SSE values. (d) Difference rhythm between consecutive SSE values.

Number of clusters	SSE	SSE diff	SSE diff rhythm
2	12,14368828	5,145574	
3	6,998114005	1,465283	3,680291
4	5,532831003	0,402046	1,063237
5	5,130784647	0,443769	-0,04172
6	4,687015166	0,929425	-0,48566
7	3,757589924	0,34968	0,579745
8	3,40790992	3,40791	-3,05823

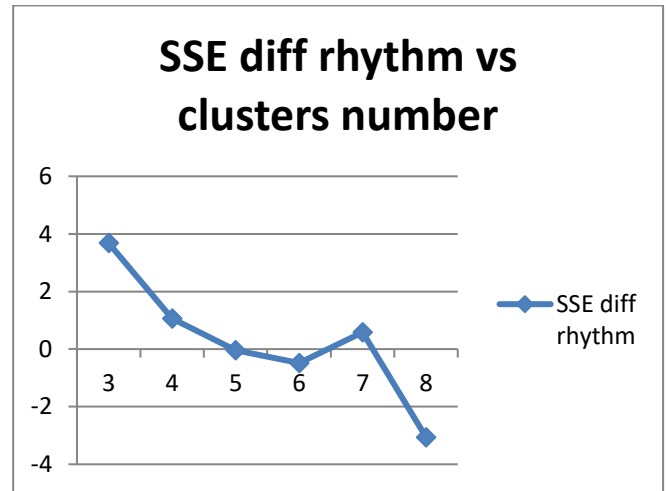


Fig.3. SEE difference rhythm vs clusters number for Iris-flower data.

We can see from the above diagram the optimal number of clusters: 3, corresponding to the higher SSE difference rhythm value.

B. Image Segmentation data.

- Creators: Vision Group, University of Massachusetts.
- Donor: Vision Group (Carla Brodley, brodley@cs.umass.edu), November 1990.
- Relevant Information: The instances were drawn randomly from a database of 7 outdoor images. The images were hand-segmented to create a classification for every pixel. Each instance is a 3x3 region.
- Number of Instances: Training data: 210 (30 instances per class). Test data: 2100 (300 instances per class).
- Number of Attributes: 19 numeric attributes and the class.
- Attribute Information: (1) region-centroid-col: the column of the center pixel of the region. (2) region-centroid-row: the row of the center pixel of the region. (3) region-pixel-count: the number of pixels in a region = 9. (4) short-line-density-5: the results of a line extractoin algorithm that counts how many lines of length 5 (any orientation) with low contrast, less than or equal to 5, go through the region. (5) short-line-density-2: same as short-line-density-5 but counts lines of high contrast, greater than 5. (6) vedge-mean: measures the contrast of horizontally adjacent pixels in the region. This attribute is used as a vertical edge detector (mean). (7) vegde-sd: (see 6, standard deviation), (8) hedge-mean: measures the contrast of vertically adjacent pixels. Used for horizontal line detection (mean). (9) hedge-sd: (see 8, standard deviation). (10) intensity-mean: the average over the region of (R + G + B)/3. (11) rawred-mean: the average over the region of the R value. (12) rawblue-mean: the average over the region of the B value. (13)

rawgreen-mean: the average over the region of the G value. (14) exred-mean: measures the excess red: $(2R - (G + B))$. (15) exblue-mean: measures the excess blue: $(2B - (G + R))$. (16) exgreen-mean: measures the excess green: $(2G - (R + B))$. (17) value-mean: 3-d nonlinear transformation of RGB. (Algorithm can be found in: Foley, James D. and Van Dam, Andries, 1982). (18) saturatoin-mean: (see 17). (19) hue-mean: (see 17).

- Class Distribution: brickface, sky, foliage, cement, window, path, grass.
- Missing Attribute Values: None

You can see below the data extracted after the consecutive K-means applications.

Table 2. Image Segmentation data. K-Means clustering in WEKA environment. (a) Number of clusters. (b) SSE: within cluster sum of squared errors. (c) Difference between consecutive SSE values. (d) Difference rhythm between consecutive SSE values.

Number of clusters	SSE	SSE diff	SSE diff rhythm
4	361	40	
5	321	31	9
6	290	27	4
7	263	14	13
8	249	10	4
9	239	21	-11
10	218	16	5
11	202	10	6
12	192	3	7
13	189	10	-7
14	179	2	8
15	177	9	-7
16	168	5	4
17	163	3	2
18	160	4	-1
19	156	3	1
20	153	1	2
21	152	8	-7
22	144	5	3
23	139	5	0

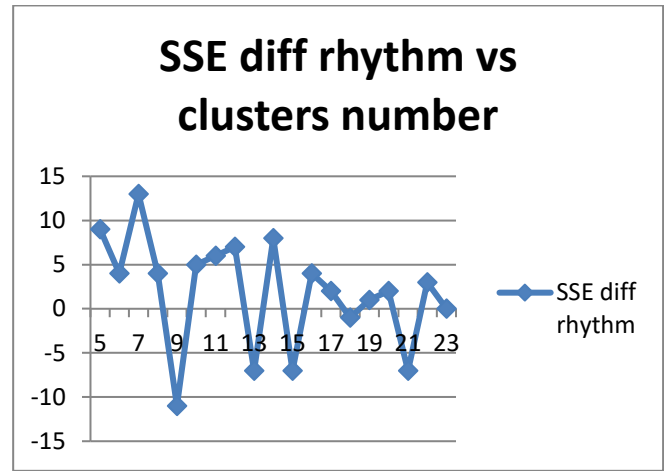


Fig.4. SEE difference rhythm vs clusters number for Image segmentation data.

We can see from the above diagram the optimal number of clusters: 7, corresponding to the higher SSE difference rhythm value.

IV. CONCLUSION

This article deals with the open question of the nature of information. Is the information inherent element of the things or just a humanly tool aiming at the better understanding of our environment? In this frame two types of information were examined. The first information type concerns the real informative part of a natural signal after eliminating the noisy components. The second information type concerns the classification of things or events, members of a wide group, into smaller groups, so that they both inherit the general characteristics and they are also, as members of the distinct subgroups, discriminated between each other. Distinction is an inherent feature of nature, while at the same time the classification process satisfies the humanly need for better understanding through specification. The outcome of the first information type is the number of harmonics of the signal, which identifies with the order of the underlying system. The outcome of the second information type is the number of clusters where a group of things is classified. These two measures seem to be related, since each of them shows the number of independent projection axes of the examined features. The innovative part of this paper is the presentation of a method for the selection of the appropriate number of clusters in the K-Means clustering process.

Capurro gives a very interesting interpretation of the concept of information through its etymological connection to the concept of form. Information in this sense is the process of forming a piece of matter or, metaphorically, human knowledge (Capurro, 1996, 256-270) [5]. Shannon and Weaver in their mathematical theory of communication establish a neutral, independent of the human comprehension, meaning of information content, that of the information entropy. Carl Friedrich von Weizsacker relates this neutral substance, which is neither matter nor energy, to the Platonic eidos and the Aristotelian form [5]. Therefore, information content could

be meant as an autonomous substance, a pattern archetype for in-forming something (Capurro, 1991) [4]. Under this view, that extracted number, either as the model order or as the multiplicity of clusters, seems to be an inherent element of a natural process, either in the frame of signals to be eliminated from noise or in the frame of a group of things to be divided through some distinctive features. Thus the harmonics of a signal or the clusters of an overall group appear as non-temporal characters of the temporal phenomena. These characters remind us of the platonic ideas, a general concept of genera and species/eide/ ideas/ forms. We could suppose that a genus itself, such as the iris flower studied in the previous section, is a part of the nature corresponding to a certain design, and holding an objective character apart from the human endeavor to categorize it as a means of better insight and understanding.

REFERENCES

- [1] Bakamidis, S.; Dendrinis, M. and Carayannis, G., "SVD Analysis by Synthesis of Harmonic Signals", *Acoustics Speech and Signal Processing (ASSP)*, IEEE, Vol. 39, No.2, pp. 472-477, Feb. 1991, <https://ieeexplore.ieee.org/document/80831>
- [2] Balasubramanian, S.; Rajavel, R.; Kar, Asuthos, "Ideal ratio mask estimation based on cochleagram for audio-visual monaural speech enhancement", *Applied Acoustics*, Volume 211, August 2023.
- [3] Carayannis, G. and Gueguen, C. "The factorial linear modelling : a Karhunen-Loeve approach to speech analysis", *Proc. IEEE Int. Conf. Acoust. Speech Signal Process., ICASSP*, pp. 489-492, 1976.
- [4] Capurro, R., "Foundations of Information Science, Review and Perspectives", *International Conference on Conceptions of Library and Information Science, University of Tampere, Tampere, Finland, 26-28 August 1991*. Also available at www.capurro.de/tampere91.htm
- [5] Capurro, R., "On the genealogy of information", *International Conference: Information. New Questions to a Multidisciplinary Concept, organized by the Chair for Philosophy of Technology at the Technical University of Cottbus held from March 1st to 3rd, 1994*. Berlin: Akademik Verlag Berlin, pp. 259-270, 1996, also available at www.capurro.de/cottinf.htm
- [6] Christensen, Knud B.; Christensen, Mads G.; Boldt Jesper B.; Gran, Fredrik, "Experimental study of generalized subspace filters for the cocktail party situation", *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 20-25 March 2016.
- [7] Dendrinis, M., "Obstructive Sleep Apnea (OSA) Detection through Energy and Spectral Measures", *International Conference on Signal Processing Applications and Technology (ICSPAT-98)*, IEEE, Toronto, Sept. 1998.
- [8] Dendrinis, M.; Bakamidis, S. and Carayannis G., "Speech Enhancement from Noise : A Regenerative Approach", *Speech Communication*, Vol.10, No.1, pp.45-57, Feb.1991 <http://www.sciencedirect.com/science/article/pii/016763939190027Q>
- [9] Deppisch, Thomas; Amengual, Sebastià V.; Calamia, Garí Paul; Ahrens, Jens "Direct and Residual Subspace Decomposition of Spatial Room Impulse Responses", *IEEE/ACM Transactions on Audio, Speech and Language Processing*, Volume 31, pp. 927–942, 2023.
- [10] Fisher, R.A. "The use of multiple measurements in taxonomic problems", *Annual Eugenics*, 7, Part II, 179-188, 1936; also in *Contributions to Mathematical Statistics*, John Wiley, NY, 1950.
- [11] Foley, James D. and Van Dam, Andries, *Fundamentals of Interactive Computer Graphics*, Addison-Wesley Publishing Company, 1982.
- [12] Rao, S. S. and Gnanapmkasam, D. C., "A criterion for identifying dominant siju1ar values in the SVD based method of harmonic retrieval," in *Proc. iCJ5SP 88* (New York, NY), E.6.5, pp. 2460- 2463, 1988.
- [13] Stainhaouer, Gregory; Bakamidis, Stelios; Dologlou, Ioannis, "Automatic Detection of Allergic Rhinitis in Patients", *International Conference on Computational Science and Computational Intelligence (CSCI)*, 5-7 Dec. 2019.
- [14] Tong, Renjie; Ye, Zhongfu, "Data fusion over localized sensor networks for parallel waveform enhancement based on 3-D tensor representations", *Signal Processing*, Volume 141, pp. 249-260, December 2017.
- [15] Wang, Fasong; Wang, Zhongyong; Li, Rui; Zhang, Linrang "An efficient algorithm for harmonic retrieval by combining blind source separation with wavelet packet decomposition", *Digital Signal Processing*, Volume 46, pp. 133-150, November 2015.