

2018

# Identifying the Sales Patterns of Online Stores with Self-organising Maps on Time Series Data

Markus Makkonen

*University of Jyväskylä, markus.v.makkonen@jyu.fi*

Lauri Frank

*University of Jyväskylä, Finland, lauri.frank@jyu.fi*

Follow this and additional works at: <https://aisel.aisnet.org/mcis2018>

## Recommended Citation

Makkonen, Markus and Frank, Lauri, "Identifying the Sales Patterns of Online Stores with Self-organising Maps on Time Series Data" (2018). *MCIS 2018 Proceedings*. 11.

<https://aisel.aisnet.org/mcis2018/11>

This material is brought to you by the Mediterranean Conference on Information Systems (MCIS) at AIS Electronic Library (AISeL). It has been accepted for inclusion in MCIS 2018 Proceedings by an authorized administrator of AIS Electronic Library (AISeL). For more information, please contact [elibrary@aisnet.org](mailto:elibrary@aisnet.org).

# IDENTIFYING THE SALES PATTERNS OF ONLINE STORES WITH SELF-ORGANISING MAPS ON TIME SERIES DATA

*Research full-length paper*

*Track: Digital Markets*

Makkonen, Markus, University of Jyväskylä, Jyväskylä, Finland, markus.v.makkonen@jyu.fi

Frank, Lauri, University of Jyväskylä, Jyväskylä, Finland, lauri.frank@jyu.fi

## Abstract

*Electronic commerce, especially in the business-to-consumer (B2C) context, has for years been a popular research topic in information systems (IS). However, the prior research on the topic has traditionally been dominated by the consumer focus instead of the business focus of online stores. For example, whereas various segmentations exist for online consumers based on their purchase behaviour, no such segmentations have been developed for online stores based on their sales patterns. In this study, our objective is to address this gap in prior research by identifying the most typical sales patterns of online stores operating in the B2C context. By using self-organising maps (SOM) to analyse the monthly sales time series collected from 399 online stores between January 2016 and December 2017, we are able to identify four approximately equally sized segments, each with its characteristic sales pattern. More specifically, two of the segments are characterised by a clear upward or downward trend in the sales, whereas the other two are characterised by strong seasonal sales variation. We also investigate the differences between the segments in terms of several key business and technical parameters of the stores as well as discuss more broadly the applicability of SOM to IS.*

*Keywords: Electronic Commerce, Business-to-Consumer (B2C), Online Stores, Sales Patterns, Self-Organising Maps (SOM), Time Series, Segmentation.*

## 1 Introduction

Electronic commerce or e-commerce, especially in the business-to-consumer (B2C) context, has for years been a popular research topic in information systems (IS). However, the prior research on the topic has traditionally been dominated by the consumer focus, such as the various factors affecting online consumer behaviour (e.g., Perea y Monsu  , Dellaert and de Ruyter, 2004; Chang, Cheung and Lai, 2005; Cheung, Chan and Limayem, 2005) as well as the segmentations of online consumers based either on self-reported survey data (e.g., Brown, Pope and Voges, 2003; Kau, Tang and Ghose, 2003; Swinyard and Smith, 2003; Bhatnagar and Ghose, 2004a, 2004b; Rohm and Swaminathan, 2004; Brengman et al., 2005; Allred, Smith and Swinyard, 2006; Barnes et al., 2007; Soopramanien and Robertson, 2007) or on click-stream data about actual visits to online stores (e.g., Moe, 2003; Su and Chen, 2015). In contrast, the business focus of online stores has traditionally gained far less attention. For example, whereas various segmentations exist for online consumers based on their purchase behaviour, no such segmentations have been developed for online stores based on their sales patterns.

In this explorative study, we aim to address this gap in prior research by identifying the most typical sales patterns of online stores operating in the B2C context. With a *sales pattern*, we refer to the temporal variation in the aggregate sales of an online store across all its sold items. Based on our review of prior literature, we find this study to be the first of its kind. Of course, sales patterns (or sales variation) as such have been investigated in several prior studies (e.g., Geurts, 1988). However, these studies have traditionally aimed at forecasting the future sales of an individual item or store instead of investigating the past sales of a wider sample of stores and identifying typicalities in them. Or those studies that done this, have done it only in the context of traditional offline retailing (e.g., Carman and Figueroa, 1986). We see this kind of an investigation and identification as highly valuable in both theoretical and practical respects. First, it helps us to better understand the characteristics of online markets and, for example, to anticipate the upcoming seasonal sales peaks, which can be considered important especially for logistic and payment service providers in terms of proper resource allocation. Second, if we are able to find causal explanations for the emergence of specific sales patterns, or at least associate them with specific business and technical parameters of the stores, we may be able to offer the stores competitive advantage by suggesting what kind of strategies help them to move away from sales patterns that are likely to be perceived as mostly negative (e.g., decreasing sales and strong seasonal sales variation) and towards sales patterns that are likely to be perceived as more positive (e.g., increasing and more stable sales).

The study is conducted in two phases by using self-organising maps (SOM) to analyse the monthly sales time series collected from 399 online stores between January 2016 and December 2017. In the first phase, we explore the different possibilities of segmenting the stores based on their sales patterns as well as the potential segment differences in terms of several key business and technical parameters of the stores. In the second phase, we concentrate on confirming the results of the previous phase.

The paper proceeds as follows. After this introductory section, Section 2 briefly discusses the basics of SOM, which has remained a relatively underutilised research method in IS. After this, we will describe the methodology and results of the study in Sections 3 and 4. The results will be discussed in more detail in Section 5, which also discusses more broadly the applicability of SOM to IS. Finally, Section 6 describes the main limitations of the study and proposes some potential paths of future research.

## 2 Self-Organising Maps

A self-organising map (SOM) or a self-organising feature map (SOFM) is a special kind of an artificial neural network developed by Kohonen (1982, 1990, 1995, 1998). SOM can be used, for example, for visualising and abstracting multi-dimensional data as well as for clustering, which generally refers to grouping objects together so that they have a maximum similarity (or minimum dissimilarity) with the objects in the same group, but a minimum similarity (or maximum dissimilarity) with the objects in the

other groups. More specifically, SOM can also be used for time series clustering, in which the grouped objects are time series consisting of multiple data points in a chronological order (Liao, 2005; Fu, 2011; Aghabozorgi, Shirkhorshidi and Wah, 2015). SOM provides a topology preserving mapping from a high-dimensional space to a low-dimensional space (Kohonen, 1995). The principle of topology preservation means that the relative distance between the mapped objects remains unchanged. In other words, the objects that are near each other in the high-dimensional space remain near each other also in the low-dimensional space. In SOM, each map consists of a network of nodes and two layers of vectors: input vectors and output vectors (Kohonen, 1995). Of these, the nodes are also commonly referred to as neurons or units, whereas the output vectors may also be referred to as weight, model, prototype, or codebook vectors. Here, we will use the terms units and codebook vectors. The input vectors are used for training the map. Thus, their number equals to the number of the objects mapped during the training. In contrast, the results of the training are reflected by the codebook vectors, each of which is associated with one unit on the map. Thus, their number equals to the number of the units to which the objects are mapped. The dimensionality of both the input vectors and the codebook vectors is the same.

Before a map is trained, its size and structure (i.e., topology) must be specified in terms of the number of units, dimensionality, dimensions, shape, and lattice (Kohonen, 1995). There is no formal rule for how many units a map should have, but the well-known and widely-used SOM Toolkit software, for example, uses an informal heuristic of  $5 \times \sqrt{N}$ , in which  $N$  is the number of the mapped objects (Vesanto, 2005). The most typical options for dimensionality is to specify the map as either one-dimensional or two-dimensional, whereas the dimensions of the map are often specified as symmetric, although also asymmetric dimensions may be used. The options for shape and lattice depend on the dimensionality of the map. In the case of a two-dimensional map, the shape of the map may be specified as a plane but also as a cylinder or a toroid, in which case either the horizontal edges or the vertical edges of the map (in the case of a cylinder) or both of them (in the case of a toroid) are joined together. Of these, toroidal maps are often preferred because they have been found to reduce the unwanted border or boundary effects that may occur at the edges of planar maps (Mount and Weaver, 2011). Respectively, in the case of a two-dimensional map, the lattice is typically specified as rectangular, in which case each unit may have neighbourhood relations up to four other units, or as hexagonal, in which case each unit may have neighbourhood relations up to six other units. Of these two, a hexagonal lattice is often preferred because it produces smoother maps that are more pleasing to the eye (Vesanto, 2005).

The algorithm that is used for training a map is an unsupervised learning algorithm, meaning that no human supervision is needed during the training (Kohonen, 1995). The training begins by initialising the codebook vectors with either random data, random samples from the input data, or samples from the input data along the linear subspace spanned by the eigenvectors with the greatest eigenvalues. After this, the input vectors are iteratively presented to the map in a random order. Each presentation triggers a process consisting of four steps. First, the distances between the presented input vector and all the codebook vectors are calculated by using the selected distance measure. The most commonly used distance measure is the traditional Euclidean distance but also other distance measures may be used. Second, based on the distances, the best-matching unit (BMU) is selected. This is the “winner” unit that has the smallest distance between its codebook vector and the presented input vector. Third, the codebook vector of the BMU is updated so that it becomes more similar to the presented input vector. The extent of updating is determined by a predefined parameter referred to as the learning rate or the learning function, the value of which decreases during the iterations in order to ensure that the training converges. Fourth, the codebook vectors of the neighbouring units to the BMU are updated so that they too become more similar to the presented input vector. The neighbourhood is determined by another parameter referred to as the neighbourhood radius or the neighbourhood function, the value of which also decreases during the iterations in order to ensure that the training converges. In the final iterations, only the BMU is typically updated. The training ends after a predefined number of iterations.

The results of the training are often visualised by using various plots (Wehrens and Buydens, 2007). For example, the codebook plot reports the resulting codebook vector of each unit, whereas the mapping

plot reports which objects were mapped to which units. In turn, the neighbour distance plot, which is also commonly referred to as the unified distance matrix or U-matrix plot, reports the sum of distances of the units to their neighbouring units. This information is often valuable because although the principle of topology preservation ensures that similar objects are mapped to the same or neighbouring units, no other plot typically reports how similar or dissimilar the units actually are to each other. The quality of the training is traditionally assessed by using two different kinds of plots. One the plots is the training plot, which reports for each iteration the mean distance between the units and their closest unit. This mean distance should decrease as the number of iterations increases, until it ultimately converges to a certain value. This value can be considered to reflect the overall homogeneity or heterogeneity of the map. In the case of no convergence, more iterations may be needed when training the map. Another plot is the quality plot, which reports the mean distances between the input vectors of the objects mapped to a unit and the codebook vector of that unit. Ideally, these mean distances should be as small as possible. Finally, the property plots provide a way to visualise the similarity or dissimilarity of the units in terms of the properties of the objects mapped to them. After training a map, the results of the training may also be further refined, for example, by clustering the units based on the similarities or dissimilarities of their codebook vectors. This can typically be done by using traditional clustering algorithms based on the distances calculated when training the map (Wehrens and Buydens, 2007).

### 3 Methodology

The data for this study was collected from 399 online stores, and it covers both their monthly sales time series from January 2016 to December 2017 as well as their key business and technical parameters at the end of January 2018. The data collection was conducted in co-operation with a Finnish electronic commerce platform provider, which uses a platform-as-a-service (PaaS) model to provide its customers an electronic commerce platform for operating the stores. The participating stores operated mainly in the Finnish market but in several product and service segments, such as clothing, jewellery and accessories, health and beauty, home and garden, sport and hobby, food and beverage, electronics, as well as automotive. All the stores had to fulfil three criteria in order to take part in the study. First, they had to be first opened before January 2015. Second, they had to be still open at the end of January 2018. Third, they had to have at least one sales transaction per month. The first two criteria were set to ensure that all the stores had been open for the whole two-year period from which the sales data was collected and that no store had opened just before or closed just after it. In other words, all the stores had more or less an equal status in terms of being relatively mature online stores, which had been in business for some time before and were not about to go out of business right after the two-year period. The third criterion was set to ensure that there was an adequate amount of sales data to be analysed from each store.

In the sales time series, the sales were measured as both sales volume (the number of sales transactions) and sales value (the value of sales transactions). Of these two, this study concentrates only on the sales measured as sales volume, but the implications of measuring the sales as sales value are addressed in the limitations. The median monthly sales of the stores ranged from a minimum of three transactions or 108 € to a maximum of 9,545 transactions or 927,476 €, with the median monthly sales of the whole sample being 34 transactions or 2,935 €. Thus, although the sample consisted mostly (about 75–80 %) of small stores with median monthly sales of less than 100 transaction or 10,000 €, it also contained very large stores in the Finnish scale. The key business and technical parameters of the stores covered their number of items, item variations, item categories, item images, campaigns, campaign items, banners, payment methods, shipping methods, customers, e-mail subscribers, and SMS subscribers. Of these, customers refer to the registered customers of the stores, while e-mail and SMS subscribers refer to the subscribers of the store newsletters via either electronic mail or short message service (SMS).

As mentioned in the introduction, the collected data was analysed in two phases by using two different-sized maps, which were specified and trained by using the *kohonen* version 3.0.4 package of R (Wehrens and Buydens, 2007). Both of them were specified as traditional two-dimensional maps with symmetric

dimensions, a toroidal shape, and a hexagonal lattice. Because no reason was seen for changing the default parameters of the package, these were used for the number of iterations, the learning rate, and the neighbourhood radius during the training (100 iterations, the learning rate decreases linearly from 0.05 in the first iteration to 0.01 in the last iteration, and the neighbourhood radius decreases linearly from a value that covers two-thirds of all the unit-to-unit distances in the first iteration to zero in the last iteration). Before the training, the raw sales time series that acted as the input vectors were standardised so that the sales of each store were measured as standard deviations from its mean sales. This is a common preparatory procedure when using SOM. In this case, it ensured that the differences in the average sales of the stores did not affect the mapping. The codebook vectors were initialised with randomly selected input vectors and the analyses were repeated with ten different random initialisations in order to ensure the stability of the results. Both the input and the codebook vectors were 24-dimensional as the sales time series covered a period of 24 months. For clustering and investigating the statistical significance of the potential cluster differences, we used the partitioning around medoids (PAM) algorithm by Kaufman and Rousseeuw (1990) as well as the Kruskal-Wallis (1952) tests and the Dunn's (1961, 1964) tests with the Bonferroni correction for multiple testing. These were run by using the *cluster* version 2.0.6 and the *dunn.test* version 1.3.5 packages of R. PAM was used as the clustering algorithm because of its applicability to a wide range of contexts, of which the context of this study was not seen as an exception. As a k-medoids algorithm, it is also more robust against outliers than k-means algorithms. In turn, the non-parametric Kruskal-Wallis tests and Dunn's tests were used because not all the assumptions for using parametric tests (e.g., normality and homoscedasticity) were met by the data.

## 4 Results

The results of the analyses are reported in the following two sub-sections, of which the first sub-section concentrates on the exploratory analysis conducted with a larger  $10 \times 10$  map and the second sub-section concentrates on the confirmatory analysis conducted with a smaller  $2 \times 2$  map.

### 4.1 Exploratory analysis with a $10 \times 10$ map

We began our analysis by training a  $10 \times 10$  map. The number of units was determined based on the number of the mapped objects ( $N = 399$ ) and the aforementioned heuristic of  $5 \times \sqrt{N}$ , which suggested that the map should have  $5 \times \sqrt{399} \approx 100$  units. After the training, we also clustered the units based on their codebook vectors by using PAM in order to highlight the potential segments. In the plots, the borders of these clusters are presented with a thicker line. The number of clusters ( $k$ ) was determined by using the overall average silhouette width (Rousseeuw, 1987) of the different clustering solutions, which evaluates both their tightness or cohesion and separation. These are reported on the left side of Figure 1. As can be seen, the four-cluster solution was found to have the highest overall average silhouette width of 0.31. The average silhouette widths of these four clusters and their sizes are reported on the right side of Figure 1. As can be seen, all the clusters were approximately equally sized.

Figure 2 presents the training plot and the quality plot, which were used to assess the quality of the training and the trained map. Here, no major issues were identified. For example, as expected, the mean distance between the units and their closest unit was found to decrease as the number of iterations increased and finally converge to about 0.028. In turn, the mean distances between the input vectors of the objects mapped to a unit and the codebook vector of that unit were found to be quite homogeneous throughout the map, with the exception of about ten units in which the mean distance was found to be somewhat higher. As expected, these units were located at the cluster borders where the neighbouring units are typically the most heterogeneous. In addition, Figure 2 presents the neighbour distance plot, which reports the sum of distances of the units to their neighbouring units. Here, about seven units were found to have a somewhat higher sum of distances. As expected, also these units were mainly located at the cluster borders, whereas the units with the lowest sum of distances were located at the cluster centres where the neighbouring units are typically the most homogeneous.

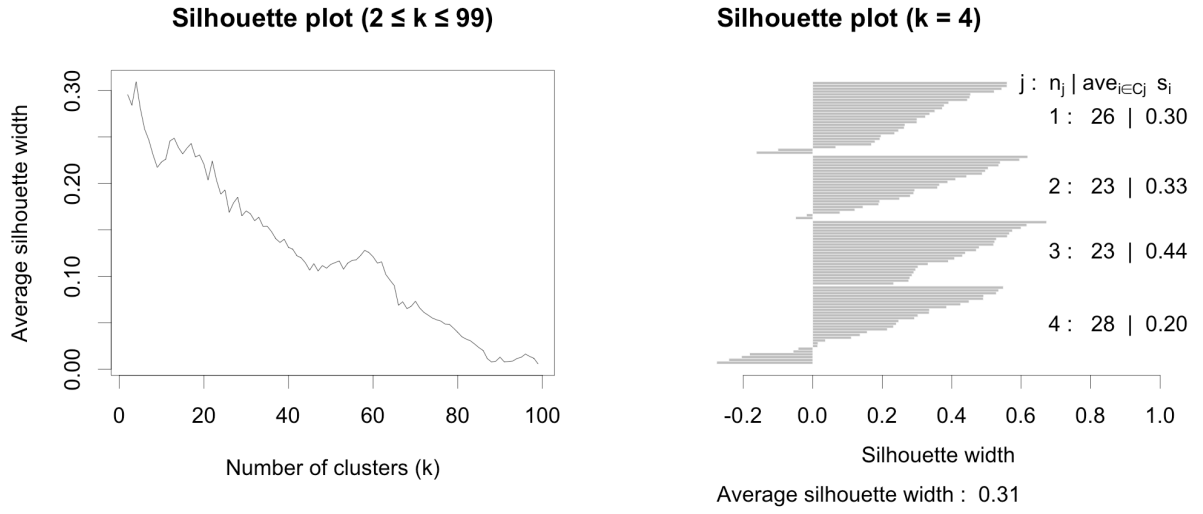


Figure 1. Silhouette plots for the  $10 \times 10$  map.

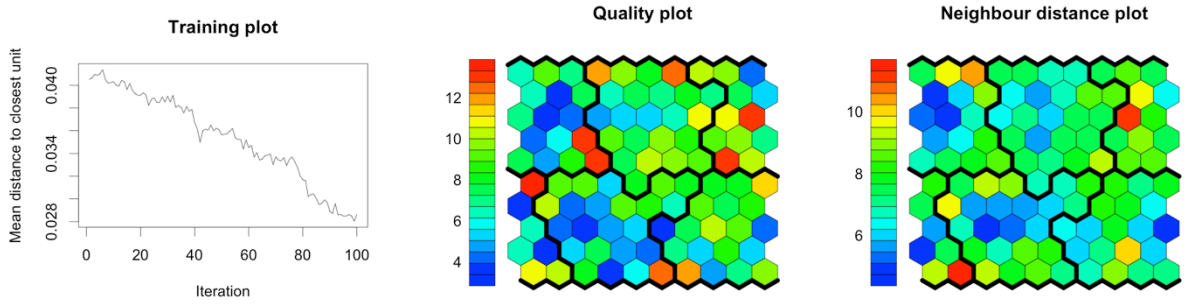


Figure 2. Training plot, quality plot, and neighbour distance plot for the  $10 \times 10$  map.

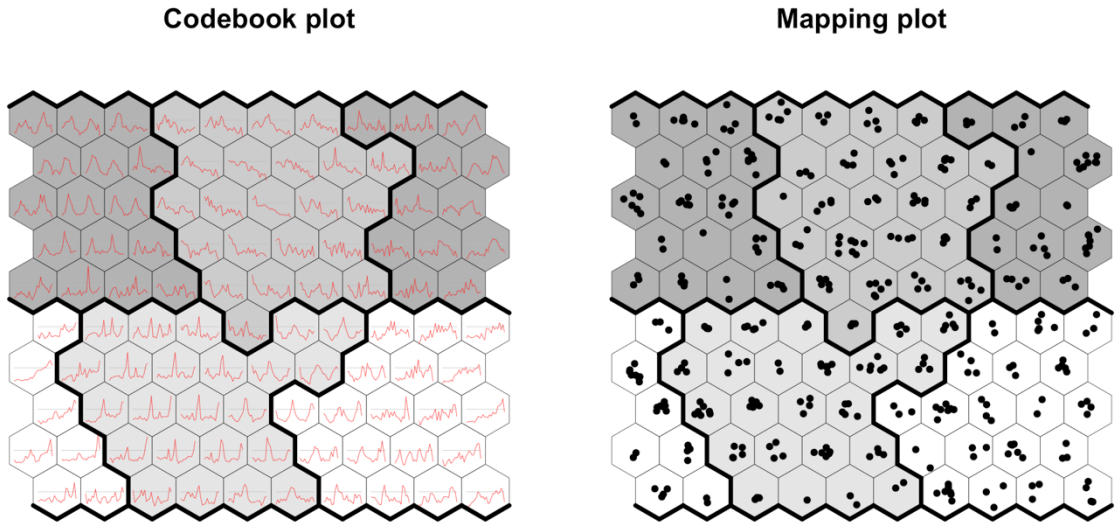


Figure 3. Codebook plot and mapping plot for the  $10 \times 10$  map.

Figure 3 presents the codebook plot and the mapping plot, which report the codebook vectors of the units and the objects mapped to them. In these plots, the clusters are highlighted also with different

background colours, with first cluster having the brightest colour and the fourth cluster having the darkest colour. As can be seen, the objects were mapped to the units relatively homogeneously throughout the map, with each unit having at least one object and no unit having more than nine objects mapped to it. We can also see that the codebook vectors of the units belonging to different clusters clearly deviate from each other in terms of shape. This is true especially for the codebook vectors of the units located at the cluster centres, which can be considered to characterise the most typical sales time series of each cluster. As can be seen, the typical sales time series in the first cluster have a clear upward trend, whereas the typical sales time series in the third cluster have a clear downward trend. In contrast, the typical sales time series in the remaining two clusters are characterised by strong seasonal variation. In the case of the second cluster, there seems to be a sharp sales peak timed at around Christmas, whereas in the case of the fourth cluster, the sales peaks seem to be less sharp and timed at the summer months. However, in both these latter cases, the sales time series lack a clear upward or downward trend.

Figure 4 presents the property plots that report the medians of 14 store parameters for the objects mapped to the units. The store parameters include the monthly sales measured as both sales volume (N) and sales value (€) as well as the 12 key business and technical parameters that were discussed in more detail above. In the case of all the parameters except for the number of payment and shipping methods, the measurement scale is logarithmic (log 10) because the medians of some units were of different order of magnitude. From the property plots, we can get a quick visual overview of the potential differences between the stores associated with each of the one hundred sales patterns. We can, for example, observe that especially the stores in the first cluster seem to differ from the stores in the other three clusters in terms of several store parameters.

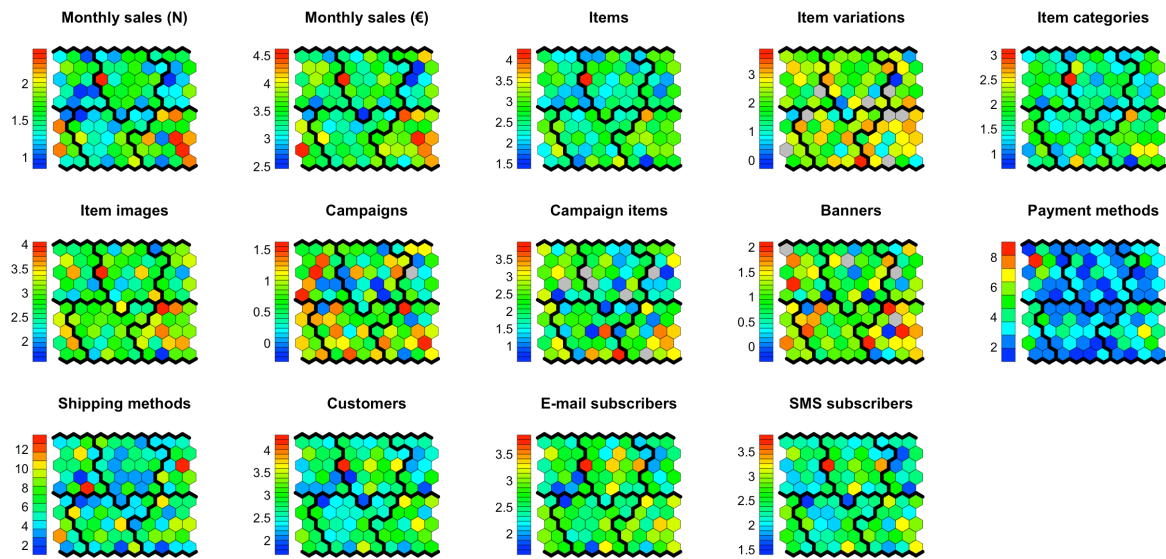


Figure 4. Property plots for the  $10 \times 10$  map.

However, the property plots do not report the precise cluster differences or their statistical significance. Therefore, two additional tables are presented. Table 1 reports the medians of the store parameters for the whole sample and for each of the four clusters as well as the differences between the sample median and the cluster medians in percentages. These support our aforementioned observation concerning the stores in the first cluster, which seem to be ahead of the stores in the other three clusters in terms of all the store parameters except for the number of payment and shipping methods.

Table 2 reports the results of the Kruskal-Wallis tests and the Dunn's tests that were used to investigate the statistical significance of the differences in the cluster medians. As can be seen, statistically significant differences between the clusters were found in all the store parameters except for the number of item variations, campaigns, and shipping methods. Most of the differences concerned the first cluster,



which was found to have higher monthly sales measured as both sales volume and sales value as well as a higher number of items and item categories in comparison to all the other clusters. The first cluster was also found to have a higher number of item images and banners in comparison to the third and fourth cluster, a higher number of customers in comparison to the second and fourth cluster, a higher number of campaign items and payment methods in comparison to the third cluster, and a higher number of e-mail and SMS subscribers in comparison to the fourth cluster. In addition, a few differences were found between the other three clusters, of which the fourth cluster was found to have lower monthly sales measured as sales volume, but higher monthly sales measured as sales value in comparison to the second and third cluster, whereas the second cluster was found to have a higher number of banners in comparison to the third cluster.

Store parameter	Median	Cluster 1 (N = 104)		Cluster 2 (N = 95)		Cluster 3 (N = 95)		Cluster 4 (N = 105)	
		Median	+ / -	Median	+ / -	Median	+ / -	Median	+ / -
Monthly sales (N)	34	64	+88 %	32	-6 %	30	-12 %	25	-26 %
Monthly sales (€)	2,935	6,424	+119 %	2,191	-25 %	2,284	-22 %	2,582	-12 %
Items	329	863	+162 %	257	-22 %	252	-23 %	272	-17 %
Item variations	167	195	+17 %	181	+8 %	136	-19 %	132	-21 %
Item categories	43	68	+58 %	36	-16 %	39	-9 %	41	-5 %
Item images	679	1,420	+109 %	617	-9 %	575	-15 %	496	-27 %
Campaigns	6	8	+33 %	6	0 %	4	-33 %	5	-17 %
Campaign items	101	210	+108 %	104	+3 %	79	-22 %	73	-28 %
Banners	10	21	+110 %	13	+30 %	6	-40 %	7	-30 %
Payment methods	3	3	0 %	2	-33 %	2	-33 %	3	0 %
Shipping methods	5	6	+20 %	5	0 %	4	-20 %	6	+20 %
Customers	494	708	+43 %	394	-20 %	537	+9 %	335	-32 %
E-mail subscribers	412	611	+48 %	416	+1 %	440	+7 %	272	-34 %
SMS subscribers	207	333	+61 %	210	+1 %	205	-1 %	151	-27 %

Table 1. Sample and cluster medians of the store parameters.

Store parameter	Kruskal-Wallis test			Dunn's tests					
	H	df	p	1 vs. 2	1 vs. 3	1 vs. 4	2 vs. 3	2 vs. 4	3 vs. 4
Monthly sales (N)	497.014	3	***	***	***	***	n.s.	***	*
Monthly sales (€)	510.778	3	***	***	***	***	n.s.	**	*
Items	17.927	3	***	*	**	**	n.s.	n.s.	n.s.
Item variations	1.848	3	n.s.	n.s.	n.s.	n.s.	n.s.	n.s.	n.s.
Item categories	13.242	3	**	*	*	*	n.s.	n.s.	n.s.
Item images	19.099	3	***	n.s.	**	***	n.s.	n.s.	n.s.
Campaigns	6.665	3	n.s.	n.s.	n.s.	n.s.	n.s.	n.s.	n.s.
Campaign items	9.736	3	*	n.s.	*	n.s.	n.s.	n.s.	n.s.
Banners	17.397	3	***	n.s.	**	*	*	n.s.	n.s.
Payment methods	9.005	3	*	n.s.	*	n.s.	n.s.	n.s.	n.s.
Shipping methods	7.584	3	n.s.	n.s.	n.s.	n.s.	n.s.	n.s.	n.s.
Customers	16.903	3	***	*	n.s.	***	n.s.	n.s.	n.s.
E-mail subscribers	17.832	3	***	n.s.	n.s.	***	n.s.	n.s.	n.s.
SMS subscribers	12.414	3	**	n.s.	n.s.	**	n.s.	n.s.	n.s.

Table 2. Results of the Kruskal-Wallis tests and the Dunn's tests for the cluster differences (\*\*\* =  $p < 0.001$ , \*\* =  $p < 0.01$ , \* =  $p < 0.05$ , n.s. = statistically not significant).

## 4.2 Confirmatory analysis with a $2 \times 2$ map

In order to confirm the aforementioned results, we replicated our analysis with a  $2 \times 2$  map, in which the four units correspond to the four clusters identified in the exploratory analysis. Figure 5 presents the training plot and the quality plot for this map. Here, the quality was found to be somewhat worse than in the case of the exploratory analysis, which was expected as the information in the 399 input vectors was now reduced to be represented by only four instead of one hundred codebook vectors. When training the map, the mean distance between the units and their closest unit was found to converge to about 0.041 instead of 0.028 but did this already in the early iterations. Respectively, in the trained map, the mean distances between the input vectors of the objects mapped to a unit and the codebook vector of that unit were found to be somewhat higher, being lowest in the unit that corresponds to the second cluster and being highest in the unit that corresponds to the third cluster of the exploratory analysis. In addition, Figure 5 presents the neighbour distance plot, which reports the sum of distances of the units to their neighbouring units. Here, the sum of distances was found to be lowest in the unit that corresponds to the first cluster and highest in the unit that corresponds to the fourth cluster of the exploratory analysis.

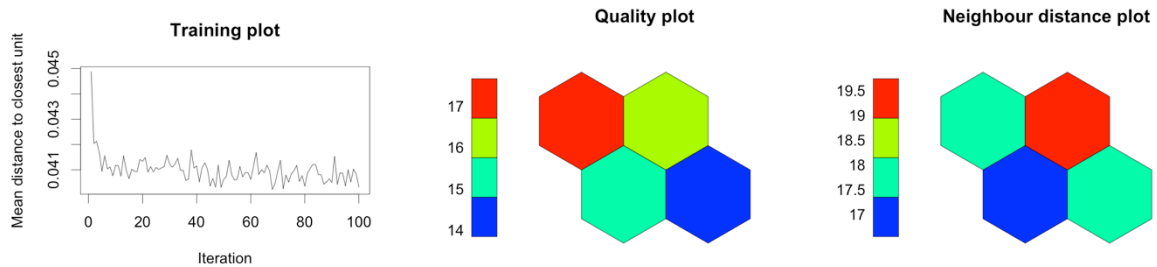


Figure 5. Training plot, quality plot, and neighbour distance plot for the  $2 \times 2$  map.

Figure 6 presents the codebook plot and the mapping plot, which report the codebook vectors of the units and the objects mapped to them. In these plots, the units are highlighted with the same background colours as the clusters of the exploratory analysis, meaning that the brightest colour represents the first unit or cluster and the darkest colour represents the fourth unit or cluster. All in all, the mapping was very consistent with the mapping of the exploratory analysis, as about 85 % of the objects were mapped to a unit that corresponds to the cluster to which the object had been mapped in the exploratory analysis. We can also see that the shape of the codebook vectors of the four units corresponds very well to the shape of the typical codebook vectors of the four clusters identified in the exploratory analysis.

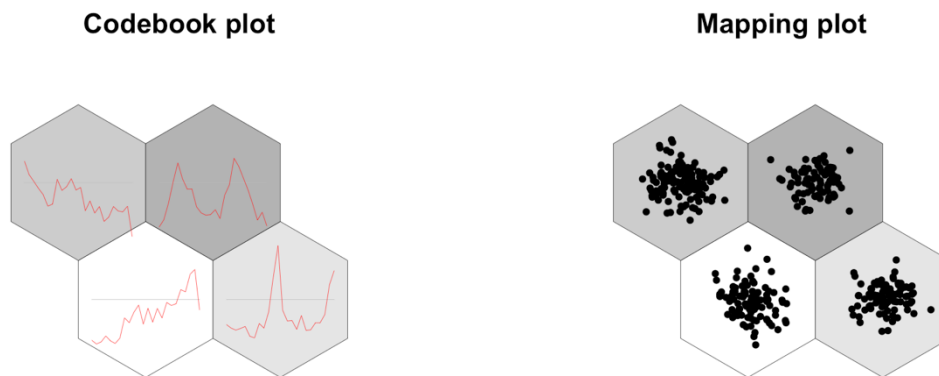


Figure 6. Codebook plot and mapping plot for the  $2 \times 2$  map.

Figure 7 presents the property plots that report the medians of 14 store parameters for the objects mapped to the units. As in the case of the exploratory analysis, the store parameters include the monthly sales measured as both sales volume (N) and sales value (€) as well as the 12 key business and technical parameters that were discussed in more detail above. From the property plots, we can once again get a quick visual overview of the potential differences between the stores associated with each of the four sales patterns. All in all, these differences are very similar to those observed in the case of the exploratory analysis, as the stores in first cluster seem to be ahead of the stores in the other three clusters in terms of all the store parameters except for the number of item variations and shipping methods.

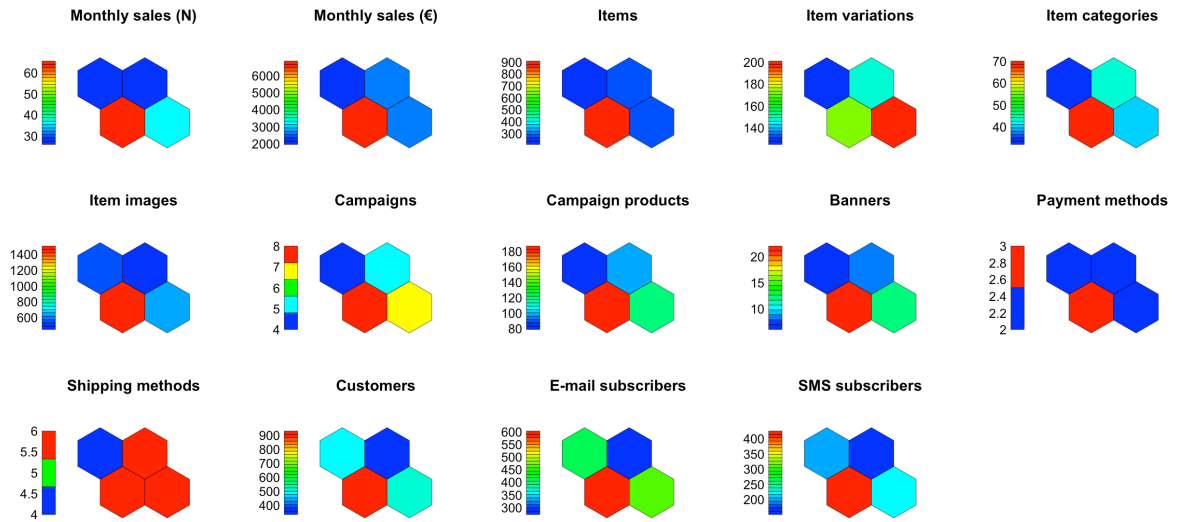


Figure 7. Property plots for the  $2 \times 2$  map.

Table 3 reports the medians of the store parameters for the whole sample and for each of the four units as well as the differences between the sample median and the unit medians in percentages. As we can see, the unit medians correspond very well to the cluster medians of the exploratory analysis.

Store parameter	Median	Unit 1 (N = 99)		Unit 2 (N = 90)		Unit 3 (N = 127)		Unit 4 (N = 83)	
		Median	+ / -	Median	+ / -	Median	+ / -	Median	+ / -
Monthly sales (N)	34	66	+94 %	37	+9 %	27	-21 %	26	-24 %
Monthly sales (€)	2,935	6,847	+133 %	2,650	-10 %	1,959	-33 %	2,706	-8 %
Items	329	908	+176 %	287	-13 %	208	-37 %	278	-16 %
Item variations	167	172	+3 %	201	+20 %	125	-25 %	149	-11 %
Item categories	43	70	+63 %	41	-5 %	32	-26 %	44	+2 %
Item images	679	1,505	+122 %	655	-4 %	575	-15 %	449	-34 %
Campaigns	6	8	+33 %	7	+17 %	4	-33 %	5	-17 %
Campaign items	101	187	+85 %	118	+17 %	79	-22 %	100	-1 %
Banners	10	22	+120 %	13	+30 %	6	-40 %	8	-20 %
Payment methods	3	3	0 %	3	0 %	2	-33 %	2	-33 %
Shipping methods	5	6	+20 %	6	+20 %	4	-20 %	6	+20 %
Customers	494	928	+88 %	502	+2 %	488	-1 %	335	-32 %
E-mail subscribers	412	603	+46 %	460	+12 %	416	+1 %	272	-34 %
SMS subscribers	207	425	+105 %	226	+9 %	195	-6 %	151	-27 %

Table 3. Sample and unit medians of the store parameters.

Table 4 reports the results of the Kruskal-Wallis tests and the Dunn's tests that were used to investigate the statistical significance of the differences in the unit medians. As can be seen, also these results are very similar to those of the exploratory analysis. However, there were a few differences. For example, statistically significant differences between the units were now found also in the number of shipping methods, which was found to be higher in the first and second unit in comparison to the third unit. In contrast, statistically significant differences could no longer be found between the first and second unit in the number of items, item categories, and customers, between the fourth and third unit in the monthly sales measured as sales volume, as well as between the fourth and second unit in the monthly sales measured as sales value. However, the second unit was now found to have higher sales measured as both sales volume and sales value in comparison to the third unit.

Store parameter	Kruskal-Wallis test			Dunn's tests					
	H	df	p	1 vs. 2	1 vs. 3	1 vs. 4	2 vs. 3	2 vs. 4	3 vs. 4
Monthly sales (N)	559.583	3	***	***	***	***	***	***	n.s.
Monthly sales (€)	652.404	3	***	***	***	***	***	n.s.	***
Items	23.954	3	***	n.s.	***	**	n.s.	n.s.	n.s.
Item variations	2.565	3	n.s.	n.s.	n.s.	n.s.	n.s.	n.s.	n.s.
Item categories	20.639	3	***	n.s.	***	*	n.s.	n.s.	n.s.
Item images	23.600	3	***	n.s.	***	***	n.s.	n.s.	n.s.
Campaigns	7.386	3	n.s.	n.s.	n.s.	n.s.	n.s.	n.s.	n.s.
Campaign items	9.255	3	*	n.s.	*	n.s.	n.s.	n.s.	n.s.
Banners	26.562	3	***	n.s.	***	*	**	n.s.	n.s.
Payment methods	9.525	3	*	n.s.	*	n.s.	n.s.	n.s.	n.s.
Shipping methods	13.584	3	**	n.s.	**	n.s.	*	n.s.	n.s.
Customers	15.462	3	**	n.s.	n.s.	***	n.s.	n.s.	n.s.
E-mail subscribers	14.421	3	**	n.s.	n.s.	**	n.s.	n.s.	n.s.
SMS subscribers	10.962	3	*	n.s.	n.s.	**	n.s.	n.s.	n.s.

Table 4. Results of the Kruskal-Wallis tests and the Dunn's tests for the unit differences  
(\*\*\* =  $p < 0.001$ , \*\* =  $p < 0.01$ , \* =  $p < 0.05$ , n.s. = statistically not significant).

## 5 Discussion and Conclusions

In this explorative study, we aimed at identifying the most typical sales patterns of online stores operating in the B2C context. By using SOM to analyse the monthly sales time series collected from 399 online stores, we were able to identify four approximately equally sized segments, each with its characteristic sales pattern. The analysis was conducted in two phases, in which the results of the first explorative analysis with a  $10 \times 10$  map were compared to the results of a second confirmatory analysis with a  $2 \times 2$  map. This latter approach has the benefit of clustering the objects directly based on their input vectors rather than indirectly based on the codebook vectors of the units, which may vary in terms of how well they represent the input vectors of the objects mapped to them. Two of the segments were characterised by a clear upward or downward trend in the sales, whereas the other two segments were characterised by strong seasonal sales variation. In one of the segments, the seasonal sales peaks lasted for about three or four months and were timed at the summer months, whereas in the other segment, they only lasted for about one or two months and were timed at around Christmas. However, most importantly, the two latter segments lacked a similar clear upward or downward trend in the sales that was found in the two former segments, meaning that the sales of the stores stayed about the same each year. Therefore, if the stores are aiming at increasing their sales in the long-term, it seems that they are likely to benefit from strategies that decrease their short-term sales variation. The explanation for this finding is most likely linked to the differences between the segments in terms of risk perceptions and investment

propensity, which are traditionally expected to be inversely related (Caballero, 1991). The stores with low seasonal sales variation obviously operate in a more certain business environment, which can be expected to have a positive effect on their investment propensity and result in increased sales if the return on investment is positive (or decreased sales if the return on investment is negative). In contrast, the stores with high seasonal sales variation obviously operate in a more uncertain business environment, which can be expected to have a negative effect on their investment propensity and hinder their growth. For example, whereas the stores with steady sales have the chance to both identify the potential changes in their economic environment and assess their return on investment throughout the year, the stores with seasonal sales can typically do this only during the on-peak months. During the off-peak months, they can often only speculate with the future while waiting and preparing for the next sales peak.

In addition, we investigated the potential differences between the online stores in each segment in terms of several store parameters. Here, especially the first segment was found to differ considerably from the other three segments. For example, the stores in the first segment were found to have higher monthly sales measured as both sales volume and sales value. Of course, this cannot be seen as particularly surprising when considering the clear upward trend in their sales. In addition, the stores in the first segment were found to be forerunners in terms of several other key business and technical parameters. These latter differences can be considered from two perspectives.

On one hand, the differences between the stores in the first segment as well as the stores in the second and fourth segment can be seen as providing insights on the business and technical parameters that are most closely associated with the reduction of seasonal sales variation. For example, the stores in the first segment were found to have a higher number of items, item categories, and item images in their selection, especially in comparison to the stores in the fourth segment. This can be expected to offer them more opportunities for diversification in terms of having in their selection also items with less seasonal sales patterns or at least items with different kinds of seasonal sales patterns. The outcomes of this diversification are likely to be seen also in the fact that the stores in the first segment were found to have more banners and registered customers as well as e-mail and SMS subscribers, once again especially in comparison to the stores in the fourth segment. After all, a more diverse item selection is likely to attract more customers to register to the store and subscribe to its newsletters. Respectively, a more diverse item selection not only offers the store more opportunities for using but also requires it to use more banner advertisement in order to gain attention for the individual items in its vast selection.

On the other hand, the differences between the stores in the first segment and the stores in the third segment can be seen as providing insights on the business and technical parameters that are most closely associated with having an upward rather than a downward trend in the sales. Many of these differences concerned the same parameters that were already discussed above. For example, the stores in the first segment were found to have more items and item categories in their selection as well as more item images and banners on their site also in comparison to the stores in the third segment. This would seem to suggest that these parameters are closely associated with not only reducing seasonal sales variation but also promoting the overall sales of the stores. This suggestion is supported by several prior studies. For example, the width and depth of the item selection have been found important especially from the perspective of capitalising the so-called “long tail” phenomenon (Anderson, 2006), which states that in many online stores, a considerable share of the sales comes not only from a few top items sold in high quantities but also from numerous niche items sold in low quantities individually, but in high quantities collectively. In turn, item images have been found important particularly in terms of promoting the flow of online shopping, the perceived usefulness of the stores, and the return intention to the stores through improved perceived diagnosticity (Jiang and Benbasat, 2004, 2007), whereas banner advertising has been found to positively affect purchase probabilities especially in the case of current customers (Manchanda et al., 2006). In addition, differences were found in the number of campaign items as well as payment and shipping methods, which were also found to be higher in the stores of the fourth segment in comparison to the stores of the first segment. These can be seen to highlight the importance of active marketing as well as offering consumers multiple alternatives in terms of payment and delivery.

In addition to the aforementioned theoretical and practical insights concerning electronic commerce and online stores, this study also presents an illustrative case example on how SOM, which has remained a relatively underutilised research method in IS, can be applied to investigate various IS phenomena. For example, from the College of Senior Scholars basket of eight IS journals, we were able to find only three studies (Lin, Chen and Nunamaker, 2000; Kiang and Kumar, 2001; Churilov et al., 2005) applying or otherwise concentrating on SOM. As we have demonstrated also in this study, the power of SOM lies especially in its ability to visualise, abstract, and cluster multi-dimensional data. One example of such data is time series data, on which we have concentrated in this study. Time series data can be considered particularly relevant for the IS context because the temporal dimension plays a central part in many IS phenomena, such as technology acceptance and use. For example, by using time series data, we can conceptualise system use not only as a simple dichotomous construct, as which it has traditionally been in conceptualised in IS theories like the technology acceptance model by Davis (1989) and the unified theory of acceptance and use of technology (UTAUT) by Venkatesh et al. (2003), but as a more complex construct capturing also the potential differences in the usage patterns of the system. This allows us to dive deeper into how factors like perceived usefulness and perceived ease of use, or the temporal changes in them, affect not only whether a system will or will not be used in the first place but also how much it will be used and what are its usage patterns. Answering these kinds of novel questions can greatly promote our present understanding on technology acceptance and use. Therefore, backed by our positive experiences from this study, we strongly encourage the usage of SOM in future IS studies.

## **6 Limitations and Future Research**

We consider this study to have three main limitations. First, the analysed sales time series covered only a two-year period and were collected only from online stores operating mainly in the Finnish market. Therefore, it is impossible to say anything about the sales patterns that exceed this two-year period or how generalisable the sales patterns found in this study are to other markets. Even the external validity of our sample to the Finnish market remains difficult to assess because although the 399 online stores included both very small and very large stores in the Finnish scale, there does not exist any reference statistics on the total number of online stores operating in the Finnish market or on their distributions in terms of sales volume and sales value. Second, we were not able to obtain from the partnering electronic commerce platform provider longitudinal data on the key business and technical parameters from the two-year period but only cross-sectional data at the end of January 2018. Therefore, we can only make associative rather than causal claims on the relationships between the store parameters and the store sales. Third, we were also not able to obtain from the partnering electronic commerce platform provider data on whether all the stores had been open for the whole two-year period or whether some of them had been temporarily closed, for example, due to maintenance. This data could have obviously been used as a valuable control variable in our study. However, although such temporary suspensions were possible, we do not believe that any of them were particularly long-lasting because all the stores participating in the study had to have at least one sales transaction per month.

Finally, it should also be noted that in addition to analysing the sales time series in which the sales were measured as sales volume, we replicated the same analysis with the sales time series in which the sales were measured as sales value. The results of this analysis were basically the same as the ones we have presented in this paper, but the identified segments were a bit more ambiguous. This was expected because the sales time series in this case were affected by not only whether a transaction was made but also by the content of the transaction in terms of how many items were purchased and for what price.

A potential path of future research would be to concentrate more precisely on the causal explanations for the sales patterns found in this study through a more thorough control of the store parameters. These parameters do not have to be strictly quantitative, as they were in this study, but they can also be more qualitative, such as assessments of the usability and aesthetics of the stores as well as information on the primary product or service segment in which the stores operate and their level of diversification.

## References

- Aghabozorgi, S., Shirkhorshidi, A. S. and Wah, T. Y. (2015). "Time-Series Clustering – A Decade Review." *Information Systems* 53, 16–38.
- Allred, C. R., Smith, S. M. and Swinyard, W. R. (2006). "E-Shopping Lovers and Fearful Conservatives: A Market Segmentation Analysis." *International Journal of Retail & Distribution Management* 34 (4/5), 308–333.
- Anderson, C. (2006). *The Long Tail: Why the Future of Business is Selling Less of More*. New York, NY: Hyperion.
- Barnes, S. J., Bauer, H. H., Neumann, M. M. and Huber, F. (2007). "Segmenting Cyberspace: A Customer Typology for the Internet." *European Journal of Marketing* 41 (1/2), 71–93.
- Bhatnagar, A. and Ghose, S. (2004a). "A Latent Class Segmentation Analysis of E-Shoppers." *Journal of Business Research* 57 (7), 758–767.
- Bhatnagar, A. and Ghose, S. (2004b). "Segmenting Consumers Based on the Benefits and Risks of Internet Shopping." *Journal of Business Research* 57 (12), 1352–1360.
- Brengman, M., Geuens, M., Weijters, B., Smith, S. M. and Swinyard, W. R. (2005). "Segmenting Internet Shoppers Based on Their Web-Usage-Related Lifestyle: A Cross-Cultural Validation." *Journal of Business Research* 58 (1), 79–88.
- Brown, M., Pope, N. and Voges, K. (2003). "Buying or Browsing? An Exploration of Shopping Orientations and Online Purchase Intention." *European Journal of Marketing* 37 (11/12), 1666–1684.
- Caballero, R. J. (1991). "On the Sign of the Investment-Uncertainty Relationship." *American Economic Review* 81 (1), 279–288.
- Carman, H. F. and Figueroa, E. E. (1986). "An Analysis of Factors Associated with Weekly Food Store Sales Variation." *Agribusiness* 2 (3), 375–390.
- Chang, M. K., Cheung, W. and Lai, V. S. (2005). "Literature Derived Reference Models for the Adoption of Online Shopping." *Information & Management* 42 (4), 543–559.
- Cheung, C. M. K., Chan, G. W. W. and Limayem, M. (2005). "A Critical Review of Online Consumer Behavior: Empirical Research." *Journal of Electronic Commerce in Organizations* 3 (4), 1–19.
- Churilov, L., Bagirov, A., Schwartz, D., Smith, K. and Dally, M. (2005). "Data Mining with Combined Use of Optimization Techniques and Self-Organizing Maps for Improving Risk Grouping Rules: Application to Prostate Cancer Patients." *Journal of Management Information Systems* 21 (4), 85–100.
- Davis, F. D. (1989). "Perceived Usefulness, Perceived Ease of Use, and User Acceptance of Information Technology." *MIS Quarterly* 13 (3), 319–340.
- Dunn, O. J. (1961). "Multiple Comparisons Among Means." *Journal of the American Statistical Association* 56 (293), 52–64.
- Dunn, O. J. (1964). "Multiple Comparisons Using Rank Sums." *Technometrics* 6 (3), 241–252.
- Fu, T.-c. (2011). "A Review on Time Series Data Mining." *Engineering Applications of Artificial Intelligence* 24 (1), 164–181.
- Geurts, M. D. (1988). "The Impact of Misrepresentative Data Patterns on Sales Forecasting Accuracy." *Journal of the Academy of Marketing Science* 16 (3–4), 88–94.
- Jiang, Z. and Benbasat, I. (2004). "Virtual Product Experience: Effects of Visual and Functional Control of Products on Perceived Diagnosticity and Flow in Electronic Shopping." *Journal of Management Information Systems* 21 (3), 111–147.
- Jiang, Z. and Benbasat, I. (2007). "Investigating the Influence of the Functional Mechanisms of Online Product Presentations." *Information Systems Research* 18 (4), 454–470.
- Kau, A. K., Tang, Y. E. and Ghose, S. (2003). "Typology of Online Shoppers." *Journal of Consumer Marketing* 20 (2), 139–156.
- Kaufman, L. and Rousseeuw, P. J. (1990). *Finding Groups in Data: An Introduction to Cluster Analysis*. Hoboken, NJ: Wiley.

- Kiang, M. Y. and Kumar, A. (2001). "An Evaluation of Self-Organizing Map Networks as a Robust Alternative to Factor Analysis in Data Mining Applications." *Information Systems Research* 12 (2), 177–194.
- Kohonen, T. (1982). "Self-Organizing Formation of Topologically Correct Feature Maps." *Biological Cybernetics* 43 (1), 59–69.
- Kohonen, T. (1990). "The Self-Organizing Map." *Proceedings of the IEEE* 78 (9), 1464–1480.
- Kohonen, T. (1995). *Self-Organizing Maps*. Heidelberg, Germany: Springer.
- Kohonen, T. (1998). "The Self-Organizing Map." *Neurocomputing* 21 (1–3), 1–6.
- Kruskal, W. H. and Wallis, W. A. (1952). "Use of Ranks in One-Criterion Variance Analysis." *Journal of the American Statistical Association* 47 (260), 583–621.
- Liao, T. W. (2005). "Clustering of Time Series Data – A Survey." *Pattern Recognition* 38 (11), 1857–1874.
- Lin, C., Chen, H. and Nunamaker, J. F. (2000). "Verifying the Proximity and Size Hypothesis for Self-Organizing Maps." *Journal of Management Information Systems* 16 (3), 57–70.
- Manchanda, P., Dubé, J.-P., Goh, K. Y. and Chintagunta, P. K. (2006). "The Effect of Banner Advertising on Internet Purchasing." *Journal of Marketing Research* 43 (1), 98–108.
- Moe, W. W. (2003). "Buying, Searching, or Browsing: Differentiating Between Online Shoppers Using In-Store Navigational Clickstream." *Journal of Consumer Psychology* 13 (1–2), 29–39.
- Mount, N. J. and Weaver, D. (2011). "Self-Organizing Maps and Boundary Effects: Quantifying the Benefits of Torus Wrapping for Mapping SOM Trajectories." *Pattern Analysis and Applications* 14 (2), 139–148.
- Perea y Monsuwé, T., Dellaert, B. G. C. and de Ruyter, K. (2004). "What Drives Consumers to Shop Online? A Literature Review." *International Journal of Service Industry Management* 15 (1), 102–121.
- Rohm, A. J. and Swaminathan, V. (2004). "A Typology of Online Shoppers Based on Shopping Motivations." *Journal of Business Research* 57 (7), 748–757.
- Rousseeuw, P. J. (1987). "Silhouettes: A Graphical Aid to the Interpretation and Validation of Cluster Analysis." *Journal of Computational and Applied Mathematics* 20, 53–65.
- Soopramanien, D. G. R. and Robertson, A. (2007). "Adoption and Usage of Online Shopping: An Empirical Analysis of the Characteristics of "Buyers", "Browsers" and "Non-Internet Shoppers"." *Journal of Retailing and Consumer Services* 14 (1), 73–82.
- Su, Q. and Chen, L. (2015). "A Method for Discovering Clusters of E-Commerce Interest Patterns Using Click-Stream Data." *Electronic Commerce Research and Applications* 14 (1), 1–13.
- Swinyard, W. R. and Smith, S. M. (2003). "Why People (Don't) Shop Online: A Lifestyle Study of the Internet Consumer." *Psychology & Marketing* 20 (7), 567–597.
- Venkatesh, V., Morris, M. G., Davis, G. B. and Davis, F. D. (2003). "User Acceptance of Information Technology: Toward a Unified View." *MIS Quarterly* 27 (3), 425–478.
- Vesanto, J. (2005). *SOM Implementation in SOM Toolbox*. URL: <http://www.cis.hut.fi/somtoolbox/documentation/somalg.shtml> (visited on 05/01/2018)
- Wehrens, R. and Buydens, L. M. C. (2007). "Self- and Super-Organizing Maps in R: The kohonen Package." *Journal of Statistical Software* 21 (5), 1–19.