

Determining the Diffusion of Innovations and the Dissemination Network in a Healthcare Organization

Dr. Duch, Jordi

Director

jordi.duch@urv.cat

Dr. Rallo, Robert

Co-director

jordi.duch@urv.cat

Dr. Guardiola, Xavier

External co-director

xavier.guardiola@simplle.cat

Vinyals, Alexandre

Student

alexandre.vinyals@estudiants.urv.cat

June 8, 2014

Contents

1	Introduction	2
2	Background Information	4
2.1	Diffusion of Innovations	4
3	State of the Art	7
3.1	Diffusion of Innovations in the Health Care Industry	8
3.2	Peer Effect Among Physicians	9
3.3	Opinion Leadership	9
3.4	Network Inferring	10
4	Data sources	12
4.1	Data correctness	15
5	Analysis	18
5.1	Identifying Diffusion Curves	19
5.2	Inferring Relationships	21
6	Results	23
6.1	Clustering methods	24
6.2	Scenario comparison	25
6.2.1	Window length	25
6.2.2	Volume of adopters	27
6.2.3	Simultaneity and connectivity correlation	28
6.3	Differentiating Opinion Leaders and Network Inferring	29
7	Conclusions	32

Chapter 1

Introduction

The goal of this research is to identify possible opinion leaders within a health care organization. To achieve this goal we use both clustering and data mining techniques, which are applied over a real world anonymized prescription database. The analysis of this prescription database provides key ingredients to infer the underlying social structure of the organization, with its possible opinion leaders. This information is inferred by analyzing both behavioral patterns and causality relationships between physicians.

Organizations are complex social entities, their social complexity can be understood with the study of large and complex social networks. Revealing information can be extracted from the link analysis of nodes. To extract those behavioral patterns and infer the causality relations, we used clustering and data-mining techniques.

The knowledge about the underlying social structure of an organization is beneficial from a business management perspective, specially if possible opinion leaders are identified. This knowledge could be used to reduce operational costs and achieve results in a more efficient way (Eg. Changing behaviors without targeting all the physicians). In fact, pharmaceutical companies who get access to this information, use it for their marketing campaigns. Possessing this knowledge, allows them to specifically target opinion leaders among the physician community. Convincing those opinion leaders to adopt their drugs increases the chances of a faster adoption rate for their drug.

New medical drugs behave as innovations, and follow diffusion processes as well. Accessing the medical prescription database allows to perform detailed analysis of those diffu-

sion curves. Applying data mining, clustering and user-profiling techniques raises privacy concerns for physicians.

This document shows a brief overview of which information could be obtained from those databases. Finally, the influence network is inferred along with the identification of possible opinion leaders. At the very end a possible solution to endorse physician privacy is discussed.

Chapter 2

Background Information

Before proceeding, this section clarifies and gives a minimum background to the reader about the concept of diffusion and how it relates with innovations.

2.1 Diffusion of Innovations

Everett M. Rogers is father of the theory *Diffusion of Innovations* [10]. The theory introduced the concept of *early adopter* and described diffusion as the way in which an innovation is communicated through certain channels over time, among the members of a social system.

Ev. Rogers closely followed diffusion processes occurring in the agricultural industry. Being born in a rural family he saw how his father loved electromechanical farm innovations, but showed reluctance for biological and chemical innovations. Thus, his father did not adopt a new corn seed which yielded 25% more crop and was resistant to drought, while his neighbor did. During the Iowa drought of 1936, while the hybrid seed corn stood tall on the neighbor farm, the crop on the Rogers farm wilted and his father was finally convinced [1].

The literal meaning of diffusion is "to spread out" and the foundation of this research is based in the analysis of the curves presented in diffusion processes. The concept of diffusion is used in a wide array of academic subjects, perhaps mostly known in physics and chemistry, where diffusion processes describe the motion of substances from high concentration areas to less concentrated areas. But the concept is applied in other disciplines as well, such as biology, sociology, finance and economics.

Ev. Rogers categorized the adopters of a new idea or product in the following categories (1) innovators (2) early-adopters (3) early-majority (4) late-majority (5) laggards . This categorization its based in standards deviation from the curve of adopters. Adoption curves are characterized as an *S-Curve* which always has a slow start and then accelerates during the intermediate phase, until finally it slows-down, getting leveraged by laggards who finally adopt the innovation. Both curves are shown in fig. 2.1, where the light grey curve shows a cumulative percentage of adopters for a given innovation, and the dark curve shows the rate at which new individuals adopt the innovation.

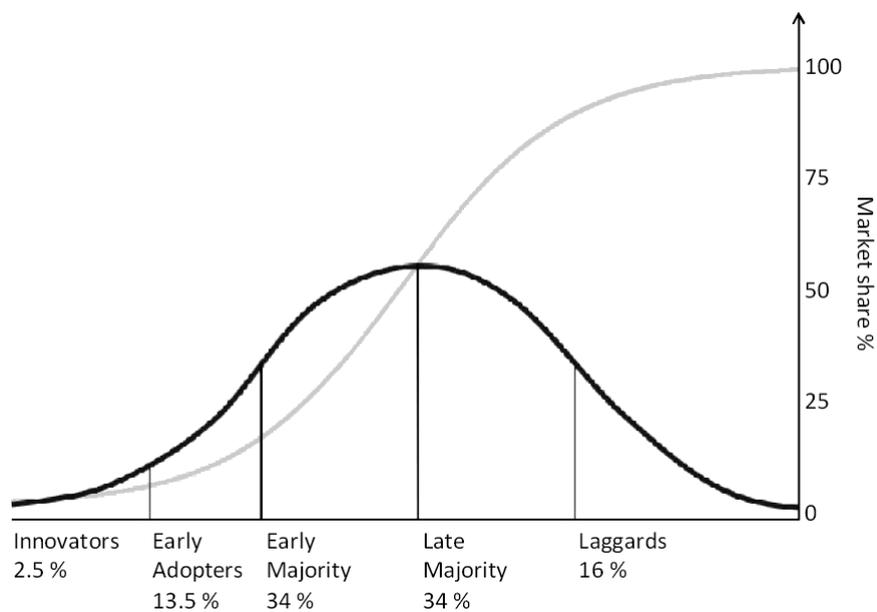


Figure 2.1: Curves of diffusion and adoption

Ev. Rogers defined that key elements in the diffusion process are (1) innovation, as the idea, practice or product that is perceived as new by an adoption unit (2) communication channels, by which messages get to one individual to another (3) social system, as a set of individuals that share a common goal or problem (4) time, as the period needed for the innovation to spread among the members of a social system .

As far as this research is concerned, those factors exist and are represented as (1) innovation, as a new medical drug (2) communication channels, as the doctor-to-doctor network, pharmaceutical marketers and medicine journals (3) social system, as the physicians working inside the health care organization (4) time, as the time required for the new medical drug to

reach a late-majority of adoption. .

If one had available a detailed description of a diffusion as it unfolds over time, in a way that one could know which individuals are under the curve at any given time, then it would be possible to infer the network by which the innovation propagated. Though all this would be under the assumption that the only communication channel for innovations to spread is the word-of-mouth.

2.2 The Role of Opinion Leaders

One of the factors that contribute to the diffusion of innovations is the influence that individuals have among others, specially since different individuals have different influences. Among those influential people, there is what its known as opinion leaders. They are found to be quite influential during the diffusion of innovations, whether it is in a positive or negative way. In fact, opinion leaders are used to accelerate diffusion processes [12] as its more efficient to target individuals who are ready to adopt innovations, than those who may not be.

Opinion leaders play a major role in the diffusion of information, knowledge or innovations among their niche of followers. For this reason, identifying opinion leaders is important from an organizational management perspective, as they can be effectively used to promote behavioral changes inside organizations [13].

Chapter 3

State of the Art

This section provides a global overview of conducted researches and discussed topics which we consider of interest for the main goal of this thesis. The scope of the thesis ranges from diffusion of innovations in health care to the privacy concerns introduced by using prescription data mining. For this reason, here we provide an overview of relevant insights found in existing academic material. This section provides evidence for

- Diffusion processes taking place in the health care industry in the form of new medical drugs
- Peering effect or influence between physicians affecting the diffusion process of new medical drugs
- Privacy concerns related with prescription data mining and physician profiling

3.1 Diffusion of Innovations in the Health Care Industry

In the health care industry, innovations are presented in the form of new treatments, surgery procedures, or new medical drugs. Successful innovations become widespread among the physician community, following a diffusion process pattern.

A recognized study called *Diffusion of Innovations Among Physicians* [4] used a combination of surveys, interviews and prescription record analysis to understand how relationships between physicians affected the diffusion of *gammanym*, a new antibiotic "wonder drug" with lesser side effects that spread out rapidly among the doctors in the medical community. This study established the importance of interpersonal networks as a communication channel for the diffusion process of the innovation.

To perform this study they used surveys and interviewed physicians to gather sociometric data. This information allowed them to directly construct the relationship network between physicians. To verify the effectiveness of the network they expected pairs of socially connected doctors to show similar behaviors. With a snowballing diffusion process — as it was occurring with *gammanym* — that is that socially connected physicians adopted *gammanym* at about the same time.

That reasoning implies that simultaneity between pairs is a consequence of social connectivity. In fact, they found that the doctor-to-doctor network operated most powerfully during the first five months of the adoption process. Doctors who still had not introduced the drug by the sixth month after its release seemed to be unresponsive to the social influence. When those doctors finally adopted it, it was because of external influences such as journals, ads or marketers... but not in response to their relationships with other doctors.

One of their conclusions was that physicians who were deeply integrated in the professional community presented a faster adoption than those who were not. And the peak effectiveness of the social links showed up during the first and second month of the adoption process, presenting a sharp decline in its effectiveness after that.

3.2 Peer Effect Among Physicians

Empirical evidence for the presence of peer effects among physicians in a health care organization has been reported in *Is There a Physician Peer Effect? Evidence from New Drug Prescriptions* [8]. This study used prescription databases from a health care organization in Taiwan specialized in patients prescribed for schizophrenia. They state that to address major challenges for empirical studies of the peer effect, their data-set had unique descriptors for patients, physicians and hospitals during a 14-year period. This allowed them to identify hospital-physician-patient pairs, which is a key factor to identify the peer effect between physicians. They found that peer effects seemed to be stronger for stable groups and that the effect intensified for larger groups, or when those groups were of the same age. This seems consistent with the observation of faster adoption rates for integrated professionals in the community.

The study reports that the measured peer effects were on maximums for newly introduced drugs, which reinforces the idea that the adoption of new drugs is driven by diffusion process. The impact of the peer effect in the context of prescription choices made by physicians has been studied as well in *Asymmetric Peer Effects in Physician Prescription Behavior: The Role of Opinion Leaders* [9], where they specifically looked for asymmetric peer effects. An asymmetric peer effect means that the influence between physicians is directed from opinion leaders to their followers, but not vice versa. They successfully quantified this effect, which manifested in the prescription patterns of physicians.

3.3 Opinion Leadership

Ev. Rogers relied in the *two-step flow of communication* model for his theory, a model that states that most people form their opinions under the influence of opinion leaders ¹.

Innovators are the kick-starters of diffusion processes, specially in the early stages of the adoption process, in which a group of innovators, opinion leaders and early adopters begin

¹The model may not be fully accurate, specially since the democratization and wide access to information available nowadays, but even with innovations reaching a wider population through mass media, and less word of mouth, persons still seek for advice from their opinion leader.

the adoption of an innovation. In fact, opinion leaders play a major role during the diffusion process, selecting ideas they would like to try and coping with the associated risks of trying new things. More significantly, opinion leaders are characterized for having plenty of social connections.

They play a major role in situations of high uncertainty, in which they build enough trust on their followers to adopt new ideas [3], contributing to the reduction of knowledge that remains unused, and making organizations progress. New medical drugs are examples of situations with high-uncertainty, providing a suitable situation to reveal the asymmetric peer effect that opinion leaders have.

In *Disseminating Innovations in Health Care* [2] they state things such as (a) find and support innovators (b) invest in early adopters (c) innovators are diamonds in the rough . They provide arguments that reinforce the importance of opinion leaders and their important role within organizations. Which somehow justifies the efforts spent in their identification.

In fact, physicians who likely are opinion leaders get targeted by pharmaceutical companies, which personally send their marketers in an effort to increase the adoption rate of their drugs, hopefully preventing the prescription of the competence drugs. The reasons they target opinion leaders are obvious, not only they become more efficient in their marketing efforts, but if an opinion leader adopts their new medical drug, chances are that physicians who follow his counsel are more likely to adopt the drug as well.

3.4 Network Inferring

The disposal of socio metric data allows to directly build the peer network. However, the underlying network in a diffusion process usually remains unobserved, as its the case in this thesis, where there is no alternative but to infer such network.

The observation and study of behavioral patterns provided key ingredients for inferring relationships in *Inferring friendship network structure* [5], where behavioral patterns from mobile phone data have been used to identify the underlying social network. As its stated that *”data collected from mobile phones have the potential to provide insight into the underlying relational dynamics of organizations, communities and, potentially, societies”*. Could the

study of behavioral patterns on a medical prescription database lead to similar results?

In *Inferring networks of diffusion and influence* [6] they state how inferring a network is possible provided that one has a detailed description about the infection of the nodes of such network. That is, a detailed description with the times in which nodes adopt pieces of information or innovations.

A detailed description about the infection of the nodes — physicians — as they get infected over time — adopt the new drug — is certainly provided in the prescription database.

3.5 Prescription Data Mining and Privacy Concerns

Prescription data mining is extensively used by pharmaceutical companies. It is said that when drug company representatives visit physicians to market their products, they already know the prescription patterns of the physician they are visiting [11]. The acquisition of this data is done by buying prescription databases from pharmacies. Medical prescriptions are filled with personal information of the patient along with the prescriber information. It is clear that applying user profiling techniques over physicians and patients raises privacy concerns for them [15].

In 2006 several states in the U.S. enacted laws to ban the use of physician prescription databases for commercial uses. Those laws were fought by large medical data collection firms until 2011, when the laws were finally struck down by the supreme court, alleging a commercial free-speech rights in violation of the First Amendment of the U.S. Constitution [14]. Therefore, physicians' privacy is still a concern nowadays. The topic has been discussed in law, medicine and ethics journals, where they suggested how those issues could be alleviated by requiring explicit consent from physicians to receiving salesman visits. It seems thought that the safety of consumers keeps being handed to an unregulated private market [7].

Chapter 4

Data sources

To achieve our goals we used a prescription database obtained from the data warehouse of a health care organization. The database contains prescription records from a population of about 200.000 citizens during a 20 months period, from 2012 to 2013. We constructed a new schema, applying extraction-transform-load techniques to the original records, the schema has the following fields (1) Anonymized patient code (2) Anonymized physician code (3) Product code (4) ATC ¹ code (5) Prescription date (6) Prescription center code (7) Patient birth yr. .

The medical drugs physicians prescribe are represented in our database with a numeric identifier. The category of this product its stored in the ATC field, which contains a classification nomenclature for medical drugs. When looking for the introduction of new products, we look at product codes appearing for the first time in some time t . A small sample of those records is provided in figure fig. 4.1, along with a numerical description of the entire data-set in fig. 4.2.

¹The ATC code has 7 characters that can be clustered in different levels *i.e the first character of the code represents the 1st level, corresponding to the anatomical main group of the drug. Three characters represent the 2nd level, corresponding to the therapeutic subgroup of the drug.*

	Patient	Physician	Product	ATC	Price	Date	Center	Birth yr.
1	ee8070d	a9c1d0a	876466	S01XA20	5,39 €	2012/01/01	01522	1931
2	66b92b4	1903b7a	831552	J01CA04	3,09 €	2012/01/01	01327	2000
3	7d52679	c55d1c7	837047	N02CA51	1,5 €	2012/01/01	00705	1950
4	43a09ca	49e5d85	815241	M01AX25	19,37 €	2012/01/01	00705	1940
5	a6525c3	f4ad493	906214	C05CA04	10,79 €	2012/01/01	00705	1952
6	742b032	a732122	788927	R05CB01	2,15 €	2012/01/01	00705	1921

...

Figure 4.1: Sample with prescription records of the database

		\bar{x}	σ	min	max
N. of prescriptions	5,192,620	-	-	-	-
N. of valid prescriptions ²	4,924,855	-	-	-	-
N. of unique products	9,310	-	-	-	-
N. of unique ATCs	966	-	-	-	-
N. of months	22	-	-	-	-
Unique patients	186,179	-	-	-	-
Unique physicians	508	-	-	-	-
Patients for physician	-	971	1,926	1	29,303
Products for ATC	-	10	21	1	196

Figure 4.2: Database descriptive attributes

The high deviation of *patients treated by physician* suggest significant differences in the volume of patients that each physician treats. Extracting the distribution of unique patients treated by physician in fig. 4.3 confirms those differences. There is a huge base of physicians whose volume of patients its 1, this could be explained by several reasons. One could be that physician students come to practice in the organization, and get to prescribe at least one time. However, this distribution confirms that there exist consolidated and integrated professionals within the organization.

²Valid prescriptions stands for prescriptions with all the required fields, *ie. prescriptions without a product code are not accepted*

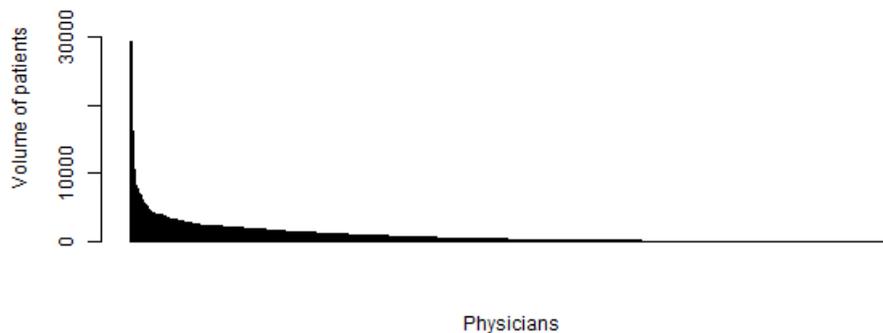


Figure 4.3: Distribution of patients treated by physicians

The same histogram is extracted to see how much products are classified in each ATC classification code, the distribution in fig. 4.4 shows significant differences again. A large amount of products in a single ATC may be related to the same drug being presented in different ways (Eg. different number of pills or different concentrations), but the reason could be a high competence in the market as well. The ATC categories with more products are for gastro-esophageal re-flux diseases, and antibacterials derived from penicillin.

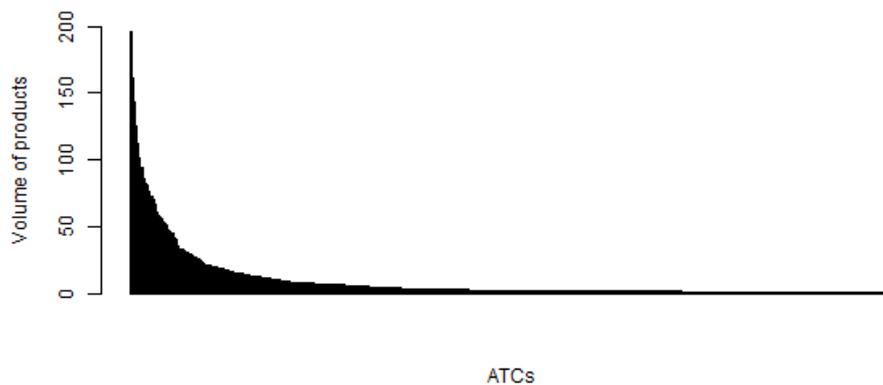


Figure 4.4: Distribution of products in ATC categories

4.1 Data correctness

This section provides a detailed description, the validity, and a better understanding of the database.

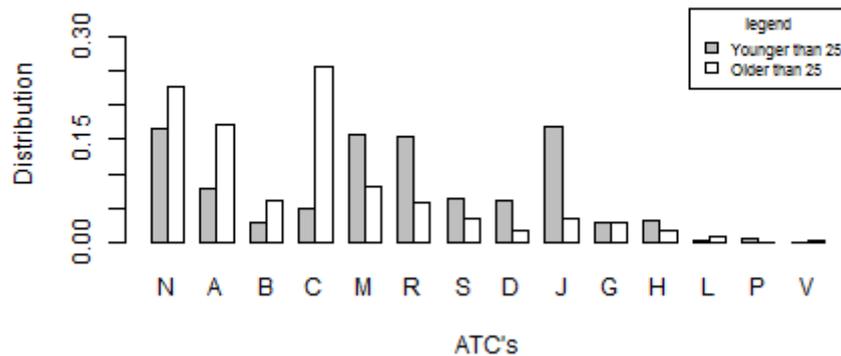


Figure 4.5: ATC distribution comparison for both groups of ages

We extracted the ATC distributions for patients, aggregating them in two major groups (a) patients under 25 years old (b) patients over 25 years old . Different distributions are expected for both groups, as patients in different ages are treated for different issues. A direct comparison of both distributions its shown in fig. 4.5.

For patients over 25 years old, (1) C - cardiovascular (2) N - nervous system (3) A - alimentary tract and metabolism are the most prescribed categories. Whereas patients under 25 years old showed more prescriptions for categories (1) M - musculo-skeletal system (2) R - respiratory system (3) J - antiinfectives for systemic use.

Cardiovascular related prescriptions represent only a 5% of prescriptions for people under 25, quite the contrary for older people, reaching a 25% in the distribution. The same phenomena is observed for the J ATC, the anti-systemic drugs category — which includes vaccines —. The comparison shows how the younger group presents a 16% proportion for J, whereas old people hardly reach a 5%. Which seems consistent, as one would expect expect people under 25 to receive more vaccines in proportion to other drugs, than the group of people over 25 years old.

In a similar way that different ages get different prescriptions, the physician specialty is expected to be reflected as well by their prescribing patterns. First hand information from the organization told us to expect two major groups of physicians (a) 70% of family physicians (b) 30% of pediatricians .

To visualize those differences we extracted the prescription patterns of all physicians, showing the distribution of ATCs they prescribe. The patterns are visualized as a heatmap in fig. 4.6, each row representing a physician, and columns representing different categories of drugs.

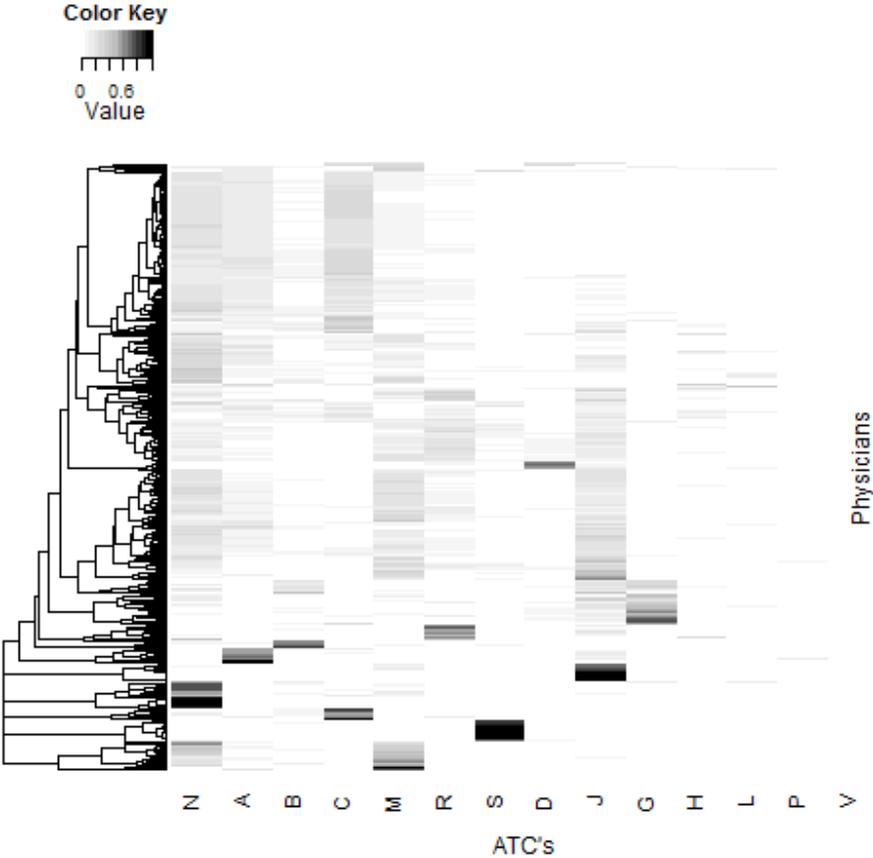


Figure 4.6: Prescription patterns for physicians

The differences in prescription patterns are clear. In fact, there exist small groups of physicians that prescribe almost entirely in one or a couple of ATC categories (Eg. physicians prescribing only in the S category, corresponding with sensory organs such as eyes and ears), they are clear examples of physician specialty at its best. The small group of physicians mostly prescribing drugs belonging to the S category probably are ophthalmologist or otologists.

After cutting the physicians tree by the fifth level, it resulted in one of the clusters containing 360 physicians, a size representing about 70.81% of the physicians in the organization, a value we were already expecting.

Chapter 5

Analysis

At this point the dataset set has been described and characterized, we looked at how products, patients and ATCs are distributed, along with the prescribing differences between physicians.

Previous literature on diffusion in health care and peer effects between physicians, reassures us that the effects are measurable and quantifiable. But perhaps more importantly, that they can be obtained by looking at prescription databases.

This thesis focuses in the inferring of relationships between physicians, by analyzing diffusion curves presented in new medical drugs. Its been shown how the peer effect accelerates the diffusion process, and how the doctor-to-doctor effects are noticeable until the fifth month. Also, socially related physicians are expected to introduce the drug at about the same time, suggesting simultaneity as a relevant factor to identify the social relationships.

The prescription database contains a complete description of any diffusion process unfolding over time ¹ during the 20 months period. This allows us to infer causality relationships between physicians, inferring how the influence cascades trough the physicians network. In order to do this, the diffusion curves for new medical drugs must be identified from the raw data.

Related literature focused in one or a few new medical drugs, with heavy support from socio-metric data that allowed them to directly observe the underlying social network of the diffusion process. Given that we do not dispose of socio-metric data, the underlying network for this study remains unobserved. In an attempt to improve the reliability of the results, we

¹As accurate as the sampling resolution permits

are going to aggregate the inferred results from individual curves, to obtain more meaningful information.

With the aggregation process, its expected that early adopters who happen to be opinion leaders, consistently adopting new drugs, will eventually get differentiated, as they will present more inferred relationships than physicians who are not, successfully identifying physicians who are likely to be opinion leaders, within the organization.

5.1 Identifying Diffusion Curves

Before proceeding, we must figure a way to identify new products and their adoption curves inside the raw dataset.

The straight forward way to identify diffusion curves is to look at the evolution of how many physicians prescribe a drug over time. Given that the data-set has prescription records from a 20 months period, to ensure that seasonality does not produce false positives (Eg. products appearing in winter due to colds would describe an adoption curve, yet they would not be an innovation or new drug) the identification will be performed for products appearing in the prescription records after the first twelve months, avoiding one-year seasonality effects.

Thus any product prescribed for the first time, a year after the first record, its a new product that could describe an adoption curve. Applying this simple restriction reduced the search scope from 9,310 to 1,103 products.

We must ensure more restrictions to hold true in order to identify diffusion curves. The problem its that those restrictions are based in *physicians over time* and the direct extraction of product histograms leaves us with *volume of prescriptions over time*. The volume of prescriptions does not represent meaningful information as far as a diffusion process its concerned, a transformation from *prescription volume* to *physician volume* has to be done. This transformation its done by extracting subsets of physicians under the curve, at each temporal step. After transforming the histograms for each product, now we can apply restrictions considering prescriber volume for new medical drugs.

A newly introduced product will follow the conditions shown in eq. (5.1). Where $A(p,t)$

represents the number of adopters prescribing the p product in time t and m_i stands for the month of introduction of the new drug, which happens to be the month in which the number of adopters prescribing p its greater than 0 for the first time.

$$\begin{cases} A(p,t) = 0, & \text{for } t \in [0, m_{i-1}] \\ A(p,t) > 0, & \text{otherwise} \end{cases} \quad (5.1)$$

For the extraction of adoption curves, we would like to control (1) granted length of the curve (2) minimum of adopters .

To grant the minimum length of n months, we check that $A(p,t) > 0$ holds true during the n previous months before the last available month in the data-set.

A minimum of k adopters its granted restricting $A(p,t) > k$ in the last available month in the data-set.

By tuning different combinations of parameters we define a set of scenarios in fig. 5.1, each one corresponded with the amount of curves that matched the conditions.

	Min. len	Cond. min. len	Min. adopters	N. curves
Scenario 1	2 mos.	$A(p, M_{20}) > 0$	$A(p, M_{22}) > 20$	58
Scenario 2	2 mos.	$A(p, M_{20}) > 0$	$A(p, M_{22}) > 10$	122
Scenario 3	2 mos.	$A(p, M_{20}) > 0$	$A(p, M_{22}) > 0$	393
Scenario 4	5 mos.	$A(p, M_{17}) > 0$	$A(p, M_{22}) > 20$	23
Scenario 5	5 mos.	$A(p, M_{17}) > 0$	$A(p, M_{22}) > 10$	42
Scenario 6	5 mos.	$A(p, M_{17}) > 0$	$A(p, M_{22}) > 0$	186

Figure 5.1: Applied restrictions and results

When we previously analyzed the prescription patterns of physicians in fig. 4.6, the smaller clusters had in average a size of 20 physicians. This makes $k > 20$ a reasonable restriction in terms of adoption widespread. Higher values would exclude specialized physicians, as a new product could not reach more physicians than themselves. The next logical restriction its no restriction at all with $k > 0$. An intermediate step is added with $k > 10$.

For the temporal length of the window, the reasonable choices are two and five months, two months being the number of months in which the doctor-to-doctor network effectiveness its at his peak, and five months as it is the period in which the doctor-to-doctor network effects are still noticeable. Examples of the identified curves are shown in fig. 5.2.

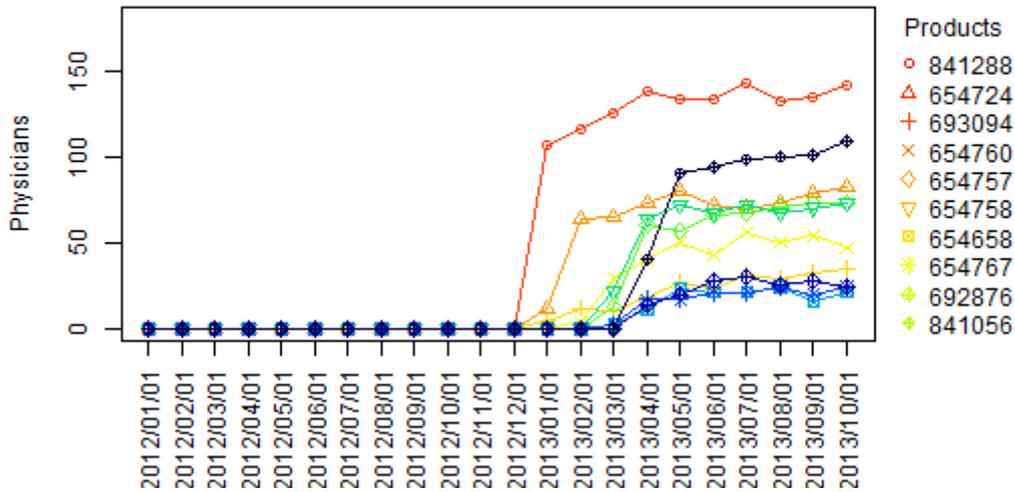


Figure 5.2: Sample adoption curves from the database

At this point one can extract the subset of physicians under the curve at any given point in time — within the one month resolution — allowing us to proceed and infer the relationships of influence between them.

5.2 Inferring Relationships

Once the diffusion curves are identified, one can produce a model to store the data of the diffusion curve as in fig. 5.3. The idea is to build a table representing the presence of the physician in each month of the curve. Rows represent the physicians under the curve, while columns the prescription volume of such physician at each month ² of the window. Physicians who join the curve past the first month are preceded by zeroes, but for reading clarity they are marked as - to clearly see when each physician joins the curve. These table is constructed for all the available curves of all scenarios.

Each temporal step has a subset of adopters who join the curve. In the invented sample physician 4 is a new adopter for the second month. Physicians 1, 2 and 3 are adopters in the first month. If we consider for this experiment that the only communication channel available for the diffusion of innovations its the doctor-to-doctor network, the only way for

²We used prescription volume to make the model more descriptive, ideally its enough with a boolean representing whether the physician prescribed or not, for each month.

	Physician	M1	M2	M3	M4	M5
1	ee3c678ee95a56a7b942dd00078b76e3	1	0	0	1	0
2	b16763ffe47918c479bc1104f0aab00d	22	43	99	91	46
3	a99e1fa6e1eee294f05d6ebf255f20b3	1	1	3	3	2
4	2dd7840f5474d7cd648c8dfb4db68a8b	-	8	2	10	2
5	278c3201e9ffacb3c1f30eac62a420d0	-	-	1	1	0
6	aaa74bef4ea7f6b7def4bc7538da45e3	-	-	1	0	0
7	f3147f883b433fd2d39c6a273e6fe911	-	-	-	1	0
8	2ddc2e49b5ebd6da78001be1560a741a	-	-	-	6	22

...

Figure 5.3: Sample of a histogram table for one specific product

physician 4 to join the curve would have been to get influenced by one — or more — of the physicians 1, 2 and 3. Basing the process in causality, physicians 1, 2 and 3 would launch arrows to physician 4 indicating an inferred connection between them. In second month, arrows would be launched from physicians 2, 3 and 4 to physicians 5 and 6 — physician 1 does not launch arrow as he is no longer present in the process by the second month —. The last arrows would be launched from the physicians under the fourth month to the new adopters in the fifth month.

The process of analyzing those tables yields both a simultaneity and adjacency matrix of physicians, examples of both shown at figs. 5.4 and 5.5.

	1	2	3	4	5	6	7	8
1	0	1	1					
2	1	0	1					
3	1	1	0					
4				0				
5					0	1		
6					1	0		
7							0	1
8							1	0

Figure 5.4: Sample simultaneity

The process may be dull for one single adoption curve, however as its been shown in fig. 5.1 we dispose of a great number of adoption curves. The aggregation of each curve matrix will eventually differentiate those physicians who are likely to be opinion leaders.

	1	2	3	4	5	6	7	8
1	0			1				
2		0		1	1	1	1	1
3			0	1	1	1	1	1
4				0	1	1	1	1
5					0		1	1
6						0	1	1
7							0	
8								0

Figure 5.5: Sample connectivity

Chapter 6

Results

This chapter presents the results obtained from the inferring process. The results are presented as heatmaps, clustered by euclidean distance. The cluster obtained from simultaneity matrix is applied to the connectivity matrix, allowing a direct comparison of both heatmaps, as dendrograms remain the same. Both the columns and rows of the heatmaps represent physicians, and the order remains the same for both. The value of each cell in the heatmaps represents the aggregated weight of the relationship, for all of the available curves.

6.1 Clustering methods

Figures 6.1 and 6.2 show the clusters formed in the simultaneity matrix, using a 5 months window with 23 curves and at least 20 adopters, in both euclidean and cosine dissimilarity distances.

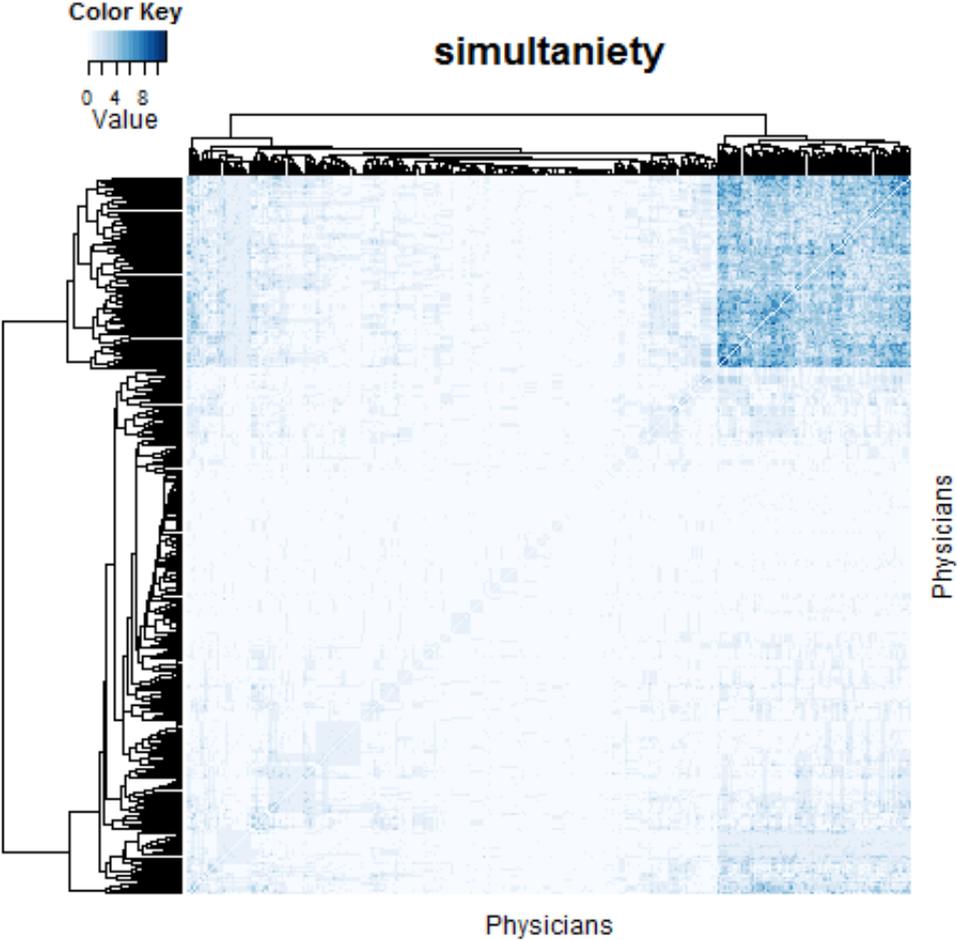


Figure 6.1: Euclidean distance

When extracting the clusters, it becomes handier to use euclidean distances, as possible opinion leaders get clustered at the highest levels of the hierarchy, whereas in cosine dissimilarity the groups are in way deeper levels, making it more tedious to handle. Also, euclidean distance provides a crisp differentiation of the clusters, which seem to fade when using cosine dissimilarity. The upcoming section will show the results using only euclidean distances.

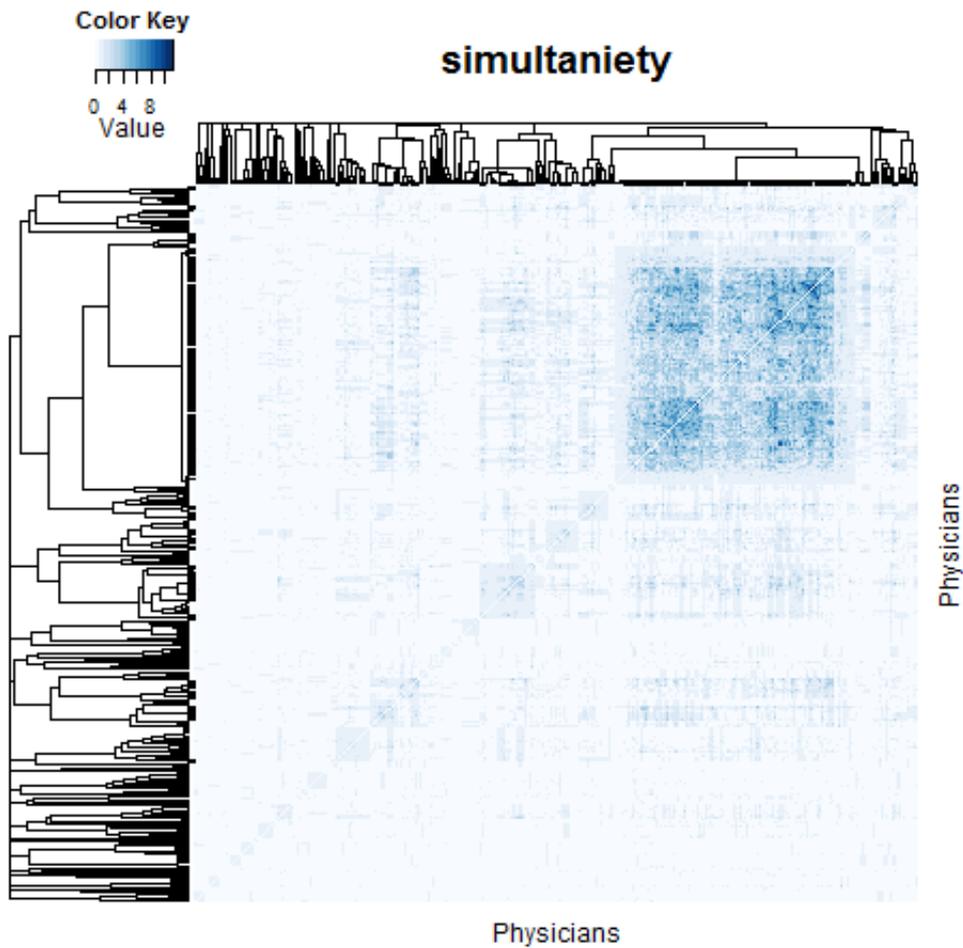


Figure 6.2: Cosine dissimilarity

6.2 Scenario comparison

6.2.1 Window length

This section shows how results vary depending on the window size. Figures 6.3 and 6.4 show simultaneity heatmaps for scenario 1 and 4 (2 and 5 months window). Figures 6.5 and 6.6 shows the respective inferred connectivity.

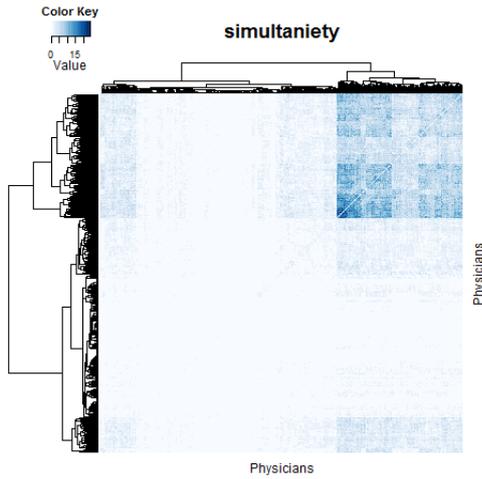


Figure 6.3: Simultaneity using 2 months window

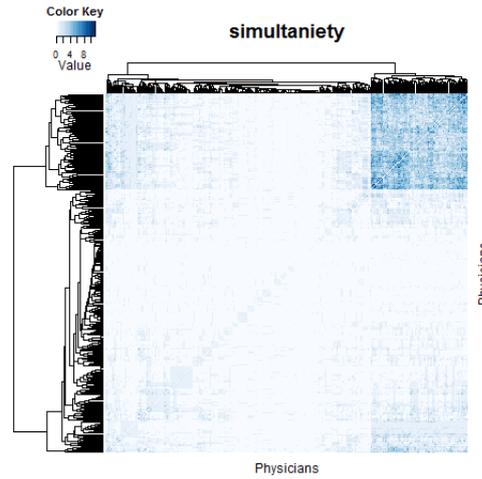


Figure 6.4: Simultaneity using a 5 months window

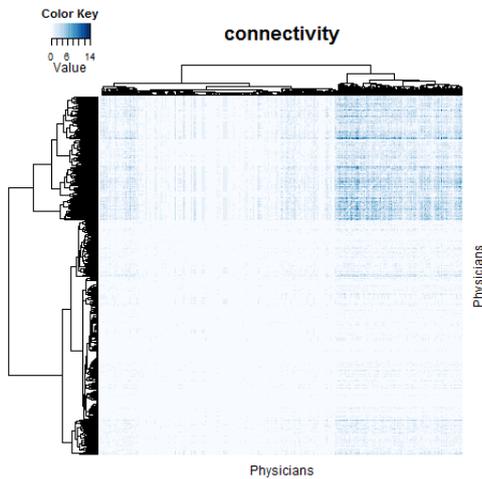


Figure 6.5: Connectivity using 2 months window

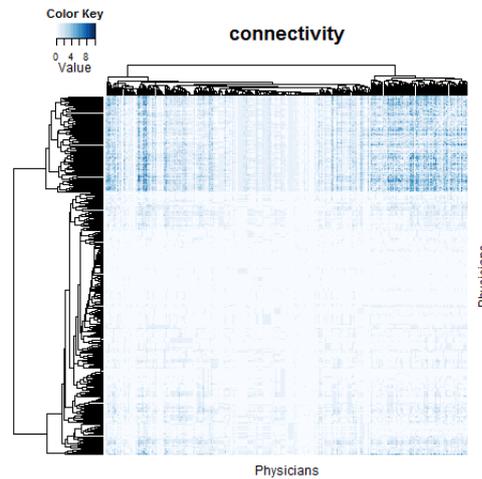


Figure 6.6: Connectivity using a 5 months window

The effects of using a wider window are notorious, inferring in a 5 months window has extra depth of inferring, which is noticeable as the connectivity reinforces horizontally, showing relations that were not inferred in the 2 months window. However, using a 2 months window allowed to better differentiate several clusters in the top-right area of the map in fig. 6.3. The reason behind this behavior is that using a 2 months window not only yields more curves (more aggregation, better inferring) but uses data in the most effective period for peering effects, obtaining more precision and less noise, unlike in the 5 month depth inferring.

6.2.2 Volume of adopters

We have seen how the length of the window affects the results. Now we present the differences that occur when curves with less than 20 adopters are used during the process. Figures 6.9 and 6.10 show the differences for a 2 months window and both $k > 20, k > 0$. Their respective simultaneity its shown below in figs. 6.7 and 6.8.

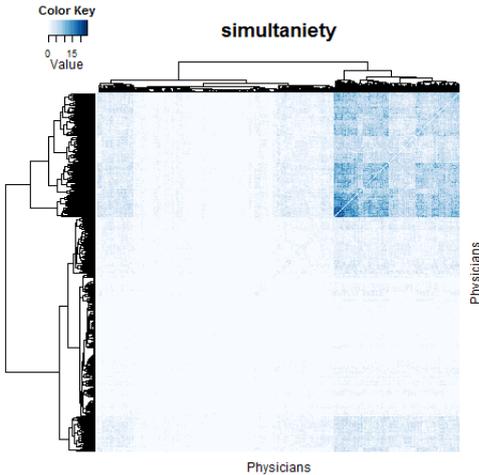


Figure 6.7: Simultaneity with > 20 adopters

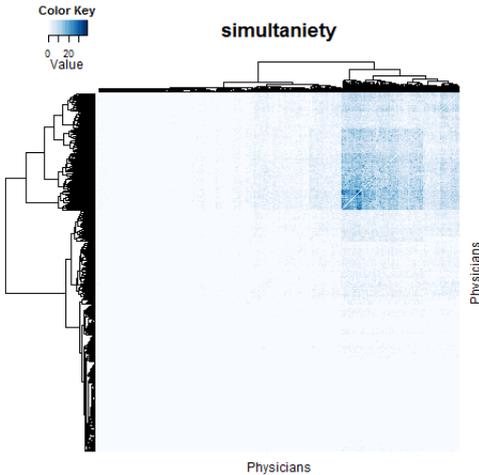


Figure 6.8: Simultaneity with > 0 adopters

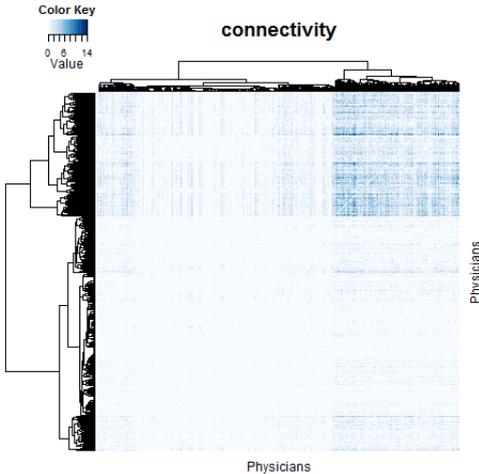


Figure 6.9: Connectivity with > 20 adopters

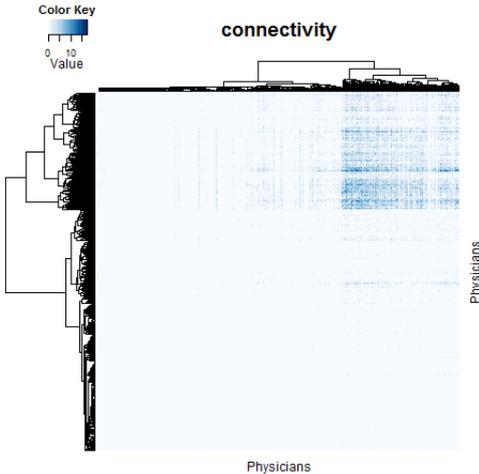


Figure 6.10: Connectivity with > 0 adopters

The implication of using $k > 0$ is that curves with small amounts of participants contribute to the aggregation process. One could think that small curves with small amounts of physicians involved would not produce significant differences, yet if those small curves involve opinion leaders, the result is what we see in fig. 6.8, that those who are likely to be opinion leaders get differentiated from the rest. Looking at the range of values, using

$k > 0$ augmented the maximum value in the scale from 25 to 40, which means that those small curves involved physicians who already had significant values of simultaneity. Opinion leaders are expected to try new things, hereby small curves with small amounts of involved physicians. Using $k > 20$ only uses curves in which the diffusion and widespread of the drug succeeds, but eliminating this restrictions does an excellent job in differentiating possible opinion leaders.

The same is seen for connectivity heatmaps, where $k > 0$ augmented the scale from 14 to 20.

6.2.3 Simultaneity and connectivity correlation

The main differences between connectivity and simultaneity are explained by (a) simultaneity is symmetric whereas connectivity loses the symmetry due to the implicit direction of the influence between physicians (b) simultaneous groups of physicians not necessarily influence others, as there could be a simultaneous group of late-adopters joining the curve. However, we can still measure the correlation between connectivity and simultaneity for the different scenarios. In figs. 6.11 and 6.12 both Pearson's and Spearman's correlation coefficients are presented for each scenario.

Scenario			Correlation summary			
mos.	adopters	curves	\bar{x}	σ	min	max
2	> 20	58	0.31	0.25	-0.19	0.77
2	> 10	122	0.34	0.26	-0.18	0.80
2	> 0	393	0.36	0.27	-0.18	0.83
5	> 20	23	0.04	0.18	-0.34	0.51
5	> 10	42	0.08	0.20	-0.32	0.60
5	> 0	186	0.12	0.22	-0.30	0.70

Figure 6.11: Pearson's correlation between connectivity and simultaneity

Previous literature told us that simultaneity was noticeable between pairs of physicians during the first and second months of the diffusion process, periods in which connected peers adopted the innovation at about the same time. Also, during the next months the effects of the doctor-to-doctor network were noticeable but not at its peak effectiveness, which weakened the simultaneity between peers.

Scenario			Correlation summary			
mos.	adopters	curves	\bar{x}	σ	min	max
2	> 20	58	0.33	0.25	-0.20	0.75
2	> 10	122	0.35	0.26	-0.19	0.76
2	> 0	393	0.37	0.26	-0.18	0.79
5	> 20	23	0.06	0.20	-0.34	0.56
5	> 10	42	0.10	0.21	-0.32	0.60
5	> 0	186	0.13	0.23	-0.30	0.62

Figure 6.12: Spearman’s correlation between connectivity and simultaneity

Similar results are manifested looking at the correlations. In the two months window there is an overall moderate relationship, with slightly strong relations appearing in maximum cases. However, with five months windows the relationships are shown to be weak overall, presenting slightly moderate relationships in maximum cases. The overall conclusion seems like 2 months window perform a better job, along with no restrictions in widespread of the drug.

6.3 Differentiating Opinion Leaders and Network Inferring

In the previous heatmaps there was one big cluster of physicians which clearly differed from the rest. This section focuses on that cluster, as it is expected to contain possible opinion leaders. As we have seen, a 2 months window with no restrictions seems to favor the differentiation of physicians, so this section focuses in scenario 3 (with its 393 curves). Applying a connectivity threshold filters about 2/3 of the total physicians. Given that the majority of physicians had none or low values for connectivity, the heatmaps become way smaller.

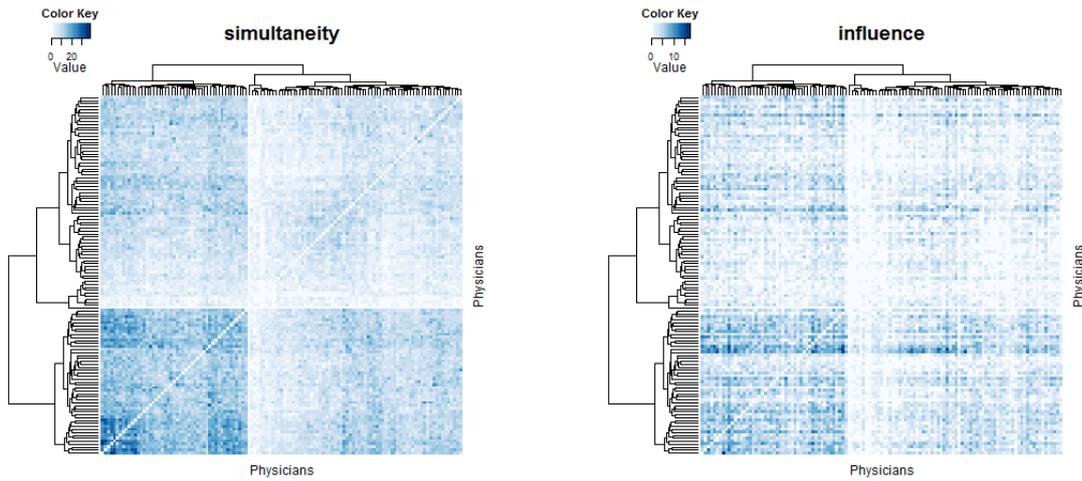


Figure 6.13: Simultaneity with > 0 adopters Figure 6.14: Connectivity with > 0 adopters

In fig. 6.13 the big cluster of physicians in bottom-left area shows how indeed simultaneity does not imply connectivity, as the same area in fig. 6.14 presents weak relationships where simultaneity was strong. Also, in the top area — where connectivity was not that strong — appeared strong lines indicating very influential physicians.

Connectivity reached around 20 for maximum cases, a reasonable threshold to obtain the possible opinion leaders could be to filter out physicians with a connectivity value smaller than 12. The filtered result is shown in fig. 6.15.

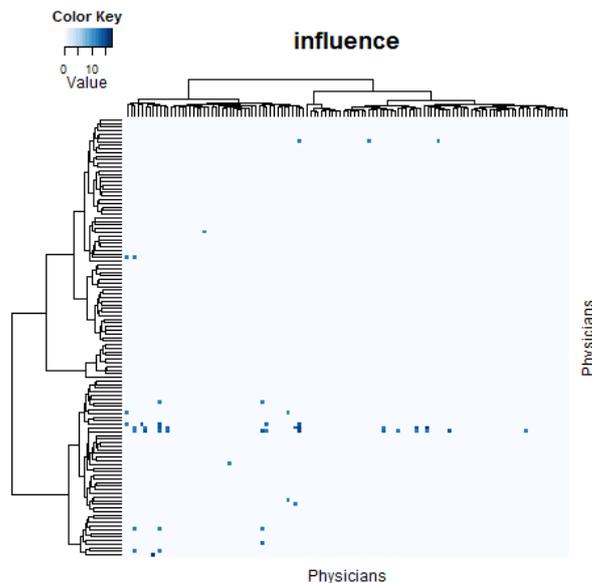


Figure 6.15: Filtered heatmap, connectivity > 12

Chapter 7

Conclusions

This thesis showed the entire process followed to infer the relationships between physicians from a health care organization, obtained exclusively from a prescription database and no socio-metric data. The results allowed the distinction of socially connected physicians which consistently adopted new drugs, which matched with the expected behaviors of an opinion leader. The inferred network also showed the overall directions of influence. All the results were anonymous as identities were previously masked behind a secure hash function. However, the organization could reveal the identities of possible opinion leaders, making this information suitable for a wide range of applications in business intelligence (Eg. one could promote or assign team leaders based on this information).

However, if external organizations such as pharmaceutical companies could get their hands in prescription databases, they could perform similar procedures with non anonymous data. Thus the privacy of physicians (and patients) could be compromised by such companies. Previous literature showed that those privacy concerns are real, and remain unsolved as of today, as pharmaceutical companies buy prescription databases from pharmacies.

Anonymizing the physician field in the prescriptions its not valid solution, as medical prescriptions are required by law to be prescriber-identified. A more complex system involving centralized government infrastructures with electronic prescriptions could effectively prevent prevent pharmaceutical companies from acquiring these databases, while still providing pharmacies ways to verify the validity of the prescription.

Europe its quite interested in deploying an electronic prescribing infrastructure, the *E-*

Prescription system. The use of electronic prescriptions has been designated as an important strategic policy to improve health care in Europe. However, due to the myriad of challenges presented in a system like that, it will be hard to see the system being fully implemented in the near-future.

Bibliography

- [1] Thomas E. Backer. Introduction. *Journal of Health Communication*, 10(4):285–288, 2005.
- [2] Donald M Berwick. Disseminating innovations in health care. *Jama*, 289(15):1969–1975, 2003.
- [3] Kenneth H Cohn and Douglas E Hough. The business of healthcare. 2008.
- [4] James S Coleman, Elihu Katz, and Hebbekt Menzel. The diffusion of an innovation among physicians. 1957.
- [5] Nathan Eagle, Alex Sandy Pentland, and David Lazer. Inferring friendship network structure by using mobile phone data. *Proceedings of the National Academy of Sciences*, 106(36):15274–15278, 2009.
- [6] Manuel Gomez-Rodriguez, Jure Leskovec, and Andreas Krause. Inferring networks of diffusion and influence. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 5(4):21, 2012.
- [7] Lawrence O Gostin. Marketing pharmaceuticals: a constitutional right to sell prescriber-identified data? *JAMA*, 307(8):787–788, 2012.
- [8] Shin-Yi Chou Muzhe Yang, Hsien-Ming Lien. Is there a physician peer effect? evidence from new drug prescriptions. 2014.

- [9] Harikesh Sasikumar Nair, Puneet Manchanda, and Tulikaa Bhatia. Asymmetric peer effects in physician prescription behavior: The role of opinion leaders. 2006.
- [10] Everett M Rogers. Diffusion of innovations. 2010.
- [11] MD Ryan M. Nunley and the Washington Health Policy Fellows. Habit-forming: Access to physician prescribing patterns, 2007.
- [12] Thomas W. Valente and Rebecca L. Davis. Accelerating the diffusion of innovations using opinion leaders. *The Annals of the American Academy of Political and Social Science*, 566(The Social Diffusion of Ideas and Things):55–67, 1999.
- [13] Thomas W Valente and Patchareeya Pumpuang. Identifying opinion leaders to promote behavior change. *Health Education & Behavior*, 34(6):881–896, 2007.
- [14] James Vicini. Supreme court strikes down state drug data-mining, 2011.
- [15] Andrew Zajac. A prescription for snooping, 2009.