

A SIMILARITY-BASED MODEL OF CONCEPT FORMATION. A PRELIMINARY ANALYSIS

Dragoş Bigu*

dragos_bigu@yahoo.com

Abstract: *In this article I undertake a preliminary analysis regarding the development of a similarity-based model of concept formation. The first part presents the new theory of concepts, based on the similarity relation. I show that only a precise concept of similarity can solve the problem of indefinite extension of the concepts formed by similarity. In the second part, I define and analyze in general lines the concept of psychological distance between objects, necessary for building a precise concept of similarity. In the third part, I show, also in general lines, how this concept can be used to build a similarity-based model of concept formation. I show the importance of “gaps” between classes at the reality level, which demarcation of categories must conform to. I briefly explain how these gaps can be captured in the model.*

Keywords: *similarity, concept, category, psychological distance between objects, gap between categories.*

According to the classical theory of concepts, their meaning is given by a set of statements that provides necessary and sufficient conditions for an object to be placed in a particular class. Since the seventh decade of the nineteenth century, this view of concepts has received a strong reply from some researchers in cognitive psychology, which showed, by the means of experimental research, that the speakers of a language do not use, in learning the concepts, necessary and sufficient criteria. Within the new theory, the similarity relation has an essential role in concept formation. Although this psychological theory has received empirical confirmation, it somehow remained at the stage of an incomplete theory.

In this article I will try to show that a similarity-based theory can be a powerful response to the classical theory of the concepts. This requires developing a model on how concepts are formed, based on the similarity relation.¹

In the limited space of this article, I can not do more than a preliminary analysis, which will help me to notice some aspects of the model, which can be detailed in a more extensive paper. In the first part, I will briefly introduce the similarity-based theory of concepts.² Then I will define and analyze the concept of

* Assistant Ph.D - *Academy of Economic Studies, Bucharest.*

¹ The discussion in this article is appropriate for both common and scientific language. However, being more precise, scientific language gives us a better environment for the formulation of appropriate examples for a theory of concepts. Therefore, I will generally use the examples in this area.

² Next, I will refer to this theory using the syntagm “*the new theory of concepts*”.

psychological distance between two objects, the key concept of a similarity-based model. In the last part, I will briefly analyze how the concept of psychological distance can be used for our purposes.

1. The new theory of concepts

The classical theory of concepts starts from the idea that any concept is characterized by a series of precise rules, which provide necessary and sufficient conditions for an object to be included in the extension of that concept. In the first version of the classical theory, that appears in its exemplary form in the works of Aristotle, these conditions took the form of a definition. Later this restriction was relaxed, and any statements containing that concept might represent such a rule. Many of the representatives of classical theory draw a clear-cut distinction between the questions regarding the meaning of concepts and psychological problems regarding the concept formation and learning. The classical theory was primarily intended to provide an answer to the first issue, detached from any psychological content. However, it could provide some clues about the psychological problems.

Regarding the problem of the meaning, classical theory must provide an answer to the question of criteria for understanding and knowing concepts. The question was when we can say that a speaker understands a certain concept. According to the classical theory, a requirement is that the speaker can determine for any possible object whether it belongs to the extension of that concept. Knowing the set of necessary and sufficient conditions that characterize the concept, the speaker will be able to identify in each case whether a particular object belongs to it.

Since the seventh decade of the nineteenth century, a number of authors have begun to question the claims of classical theory. They showed the important role of similarity between objects that belong to the same concept, completely overlooked in the classic theory. Supporters of the new theory of concepts are not in the same way as those of classical theory, adepts of a clear-cut distinction between the problem of meaning and the psychological problem of learning concepts. Therefore, they are often quite reluctant to use the word “meaning”, which leads to the idea of existence of some problems completely different from the psychological problems of concept learning and formation. Also in contrast with the followers of classical theory, the supporters of the new theory of concepts start from the real objects that belong to a certain concept. Language does not work by placing all possible objects in some categories, but starting from real objects.

The similarity relation, as used by these authors, should be seen as primitive, i.e. the relation “x is similar to y” will not be seen as an abbreviated form of the relation “x is similar to y regarding properties $P_1, P_2, P_3 \dots$ ”. If similarity concerned some properties, this would mean that the language would be built based on some relevant criteria of similarity. This would reduce the new theory to the classical one.

One of the leading critics of classical theory was Ludwig Wittgenstein. He showed that objects do not belong to the same concept in virtue of an essence they share, but in virtue of some “family resemblances” between them. Subsequently,

based on these suggestions, some cognitive psychologists have developed a theory of concepts based on the similarity relation.

Eleanor Rosch and Lawrence Barsalou have an important place among them. Rosch's research starts from the empirical observation that speakers of a language consider that some of the members of a certain category¹ are better representatives of that category than others. So, between members of a class there is a gradual transition from those who belong in the greatest degree to that class to those placed near boundaries. But if a concept were defined by a set of necessary and sufficient conditions, then any object that would satisfy those conditions would belong to that concept, and the phenomenon of gradual transition could not be explained. The conclusion of researchers is that within each category there is a prototype, an object that instantiates paradigmatically that category, and the other members of that category resemble in different degrees with it.²

In the field of philosophy of science, Thomas Kuhn shows, during the same period, that the similarity relation plays a key role for scientific concepts. Thus, in most cases, scientists do not learn to use scientific concepts by using a set of precise rules, but by using the similarity relations between objects, situations, etc.

Although the empirical observations that support the new theory of concepts are difficult to challenge, this theory is not often considered a real alternative to the classical theory. There are several rationales for this, and now I will refer to the most important of these.

This theory seems to have difficulties in explaining concept formation. To some extent, any two objects resemble each other, which will make the concepts indefinitely extend, a precise demarcation of concepts being impossible. This problem is sometimes called "the problem of wide-open texture".³ If we want to set a "boundary", it would be impossible, as the concepts defined in this way will be vague. This poses no serious difficulties when it comes to common language concepts, which do not require strict rules for use. But when it comes to scientific concepts, or other concepts used in contexts where precision is a necessary requirement, for instance concepts used in regulations, the similarity relation is not sufficient.

The main condition to reply to this criticism is to develop a similarity-based model of concept formation. The first step is to build a precise concept of similarity. I will try to do this in the next part.

2. The concept of psychological distance between two objects

On what we rely when we say that two objects are similar? Within the new theory of concepts, the answer can not be based on a set of properties that act as necessary and sufficient conditions, indicating that two objects are similar. However, clarification of the concept of similarity between two objects must rely,

¹ In the following, I will use this term, "category", to designate a class of objects that correspond to a concept. This term is used in cognitive psychology and is close to that of "natural family", as used by Thomas Kuhn.

² For a more detailed presentation of theory of Rosch and his followers, see Andersen et al., *The Cognitive Structure of Scientific Revolutions*, pp. 9-12.

³ H. Andersen, „Kuhn's Account of Family Resemblance: A Solution to the Problem of Wide-Open Texture”.

in one way or another, on the properties of those objects. One method is to use the concept of “psychological distance between two objects”. This concept, developed in the field of cognitive psychology, shows the degree of similarity between two objects: the smaller the psychological distance is, the higher the degree of similarity between objects is. Psychological distance between two objects is defined, analogue to the Euclidean distance in n-dimensional space, as follows: $d(o_1, o_2)^2 = \sum w_i (x_{1i} - x_{2i})^2$, where x_{1i} and x_{2i} are some values that characterize the objects o_1 and o_2 , from the point of view of certain characteristics, and w_i are some values representing the weight (importance) given to these characteristics. For simplicity, the weights will be “normalized”, i.e. their sum will be 1.¹

The general idea is to use this concept to identify similar objects, among those belonging to a wider class. They will be placed in the same category, while two different objects will be placed in different categories. In the last part, I will discuss about the general lines of a procedure for this. Until then I will analyze the formula for the distance between two objects.

First I will give an example. Suppose we want to establish how similar two birds are. In order to do this, we should consider a number of characteristics. Some of them will be represented by continuous variables (which can take any value) for example the size of the bird. Other ones will be represented by the cardinal variables (which can take only integer values), for instance the number of vertebrae. A third category of characteristics will be represented by categorical variables (which can take only two values, true and false), for example the ability to fly.

The above formula contains two components: the values that characterize certain aspects of objects and the weights ascribed to them, depending on their importance. In the following, I will refer to them. The first problem that arises regards characteristics that should be taken into account in measuring the psychological distance. They depend on the aims for which the demarcation of concepts is realized. This does not mean that the choice will be completely arbitrary. There are independent grounds to take into account certain characteristics. For scientific concepts, the most important reason is the occurrence of those characteristics in some scientific laws.

The problem becomes more complex when one considers that the same class of objects can be divided in several ways, based on different characteristics. For example, the class of chemical elements can be divided into metals, non-metals and semi-metals, based on a certain set of characteristics, but also in elements in main groups and those in secondary groups, based on another set of characteristics. This raises a further problem, the relationship between characteristics used together for a particular division.

On the one hand, they must not be logically deductible among them. If they were so, some of the characteristics would have a weight artificially increased. Suppose, for example, that in the example above regarding the similarity of the birds, we would take into account a cardinal variable for the number of vertebrae and also a categorical variable for the characteristic of having more than 20

¹ Formula as such is widely used in the field of cognitive psychology. (See for example G. Murphy, *Big Book of Concepts*, p. 67.). The discussion that follows is my own contribution.

vertebrae. In this case, the identical number of vertebrae would appear twice in the formula, and therefore it would be artificially overvalued. On the other hand, these relations must not be incompatible, because in this way it would be artificially undervalued. However, these observations provide only part of the answer. A more elaborate answer requires a more detailed analysis, which I will not undertake now.

A second element to be discussed is the way in which, for each characteristic, we obtain numerical values of the formula. Also in order not to overvalue certain characteristics in the formula, the values of these variables will be mapped to the same interval, for example the interval (0,1), and transformed into dimensionless values (no units of measure), in order that the distance formula be meaningful.

However, these restrictions do not lead to a single solution. At least in the case of continuous variables, a solution would be statistical normalization. Statistical normalization is an operation by which a continuous variable on a given interval is transformed into another, on a different interval, in our case (0,1), using a linear function (of the type $ax + b$). The restriction that the function is linear has a major rationale. This can provide “neutrality” of the choice of function. According to this solution, if, for instance, we have to solve the problem of values corresponding to the size of a bird, the solution is to assign value 0 to the smallest bird and value 1 to the largest bird. Values for sizes of the other birds will result from here on the basis of the restriction that the transformation function is linear.

However, this simple solution is not adequate. The main reason is that in some cases the normalization does not lead to the desired result from an intuitive point of view. For example, if we choose the solution of normalization, small differences of color, even those under the threshold of perception, would be taken into account. However, at least in some situations, this is inadequate, precisely because that difference is imperceptible. Also other examples of this type can be given, showing that the restriction of linearity of the function should be relaxed. Within this paper, I can not perform a more detailed analysis of the way in which this can be done.

Weights, the second component of the model, are necessary because some characteristics are more important than others for the appraisal of similarity degree. These weights will show the greater importance of some characteristics. For example, in the current biological taxonomy, color is considered a less important characteristic than the presence or absence of vertebral column.

Before seeing, in general lines, how the concept of psychological distance can be used, I must take some precautions regarding the hopes that I have from using this concept. First, I do not pretend that the people who use the concept of similarity and the concepts based on this use effectively this formula. Rather, what I do here is a reconstruction of this concept. Secondly, it is difficult to give to this concept a cardinal significance, i.e. to give significance to the actual value of the distance. Rather, it has an ordinal significance: a distance is greater than another, which means that two objects are more similar than the other two. Thirdly, even this use makes sense only in a determined context, in a common universe of discourse. We can say that two cats are more similar one to each other than other two cats are one to each other. But we can not say that two cats are more similar

one to each other than two seats are one to each other. The rationale for this is that in the first use, the characteristics used in the evaluation of distance are the same. But the characteristics of seats, relevant in the calculation of distance, are very different from those of cats.

3. The role of the gap between concepts in a similarity-based model

Using the above formula, many models concerning the demarcation of concepts can be built. For this reason we need a set of criteria against which to evaluate these models. Next, I will argue for a condition that has to be met by these models.

The demarcation of concepts has to conform to the “gaps” at the reality level. In order to justify this condition, a further discussion is needed. In nature, objects appear in discrete groups rather than in a continuous range. This means that there are a number of objects with a high degree of similarity among them, each of which being similar in a much smaller degree to other objects. This provides a justification for a classification based on the similarity relation. If there were no such “gaps” between categories, a demarcation of categories would be possible, but it would have no justification, since it would be only the result of setting some arbitrary “boundaries”¹. Below, I offer an example regarding the importance of this gap.

The category of acids consists of similar chemical substances, different enough from the category of bases. This will guide the research in the sense that scientists will focus on a number of properties that distinguish the two classes of compounds. Because of significant differences between these classes of compounds, more different properties are likely to be discovered. Suppose that some substances, placed in terms of properties between these two classes, are discovered. In this case, the gap between the two categories will be “covered”, which will make the research more difficult. Scientists will not know which properties, those of acids or those of bases, have to be extended to this new class. Also, even among the new substances, some will be more similar to acids and other more similar to bases. Thus, formation of a new class will have no methodological value, because scientists will not have any expectations regarding the extension of the characteristics of acids and bases. The argument above has the same relevance and for the common language. Also there, just as in scientific language, a category involves a set of expectations, and this is the justification for the demarcation of the categories.

The gap between categories can solve the problem of indefinite extension of the concepts, imposing a limit on them.² In the absence of such gap, the similarity could not lead to the formation of concepts, because the concept would indefinitely extend.

At this point a remark is necessary. Above, I have talked about the gaps between concepts as existing at the reality level. In what sense can we talk here about “reality”? Formula for psychological distance between objects can help us to clarify the issue. This depends on two elements. The first of these refers to the real

¹ T. Kuhn, *The Structure of Scientific Revolutions*, p. 45; Kuhn, „Second Thoughts on Paradigms”, pp. 312-313.

² H. Andersen, *op.cit.*

characteristics of objects and the second to which one is taken into account and with what weights. The first element depends on the way in which things are, while the second depends, at least to a certain extent, on the aims of individual (community) who demarcates the concepts. Keeping this explanation in mind, we can talk about “the reality level”, but not in a strong sense.¹

The example and the argument above show two things. First, between objects belonging to different categories must be a gap. Secondly, this gap must be respected in at the language level. So, regarding the gap between concepts, one must distinguish between two problems. The first of these problems is to develop a procedure by which gap between concepts, as they appear in measuring distances between objects, to be respected at the language level. The second issue is how large should be the gap for the demarcation of classes to have a justification.

Regarding the first issue, the aim will be to capture natural gaps, regardless of their size. The solution is to develop a method for placement of objects, by using the formula for the psychological distance between objects. Without being able to prove it now, I think an appropriate method is the following one. Each of the objects will be introduced in the category that is closer to. In a class so defined, the distance between two objects can be no matter how large, as long as the space between them is “covered” by a number of other objects that “connect” the two objects. This method should be developed in detail.

Regarding the second issue, the question regards the size of the gap between categories. We can say that a gap is a significant difference between the distances between members of a category and the distances between members of two distinct categories. The problem, which I will not discuss in this paper, will be to determine more precisely the meaning of the term “significant.” Solving those two problems will lead to a complete similarity-based model.

The main purpose of this paper has been to draw the general lines of a similarity based model of concept formation. The analysis in this paper gives me hope that such a model can be built. In this way, the new theory of concepts can be a serious competitor for the classical theory.

REFERENCES

1. Andersen, H., (2000) „Kuhn’s Account of Family Resemblance: A Solution to the Problem of Wide-Open Texture”, *Erkenntnis*, vol. 52, nr. 3, pp. 313-337.
2. Andersen, H., Barker, P., Chen, X., (2006), *The Cognitive Structure of Scientific Revolutions*, Cambridge, Cambridge University Press.
3. Kuhn, T., (1977) „Second Thoughts on Paradigms”, in *The Essential Tension*, Chicago, University of Chicago Press, pp. 293-319.
4. Kuhn, T., (1996), *The Structure of Scientific Revolutions*, Chicago, University of Chicago Press.
5. Murphy, G., (2002), *Big Book of Concepts*, MIT Press.

¹ Using different sets of characteristics can lead to different categorizations, See Andersen et al., *op. cit.*, p. 27.