# Obtaining the consensus and inconsistency among a set of qualitative facts

*Adolfo Guzman-Arenas, Adriana Jimenez-Contreras*
a**.**guzman@acm**.**org, *dyidyia@yahoo.com*

*Centro de Investigación en Computación (CIC), Instituto Politécnico Nacional (IPN), México*

*ABSTRACT.* It is well understood how to compute the average or centroid of a set of numeric values or vectors, as well as their variance. In this way we handle inconsistent measurements yielding several numeric measures. We wish to solve the analogous problem on qualitative data: How to compute the "average" or consensus of a set of non-numeric facts or observations?

This paper provides a method, based in the theory of confusion, to assess the inconsistency among a set (a bag, in fact) of qualitative observations (as opposed to quantitative measurements). Also, the most plausible value or "consensus" is determined. More than one "consensus" is at times possible. The most conspicuous outlier is determined, too.

Our approach differs from classical logic in that this logic considers inconsistency to be a number between 0 and 1; from the Theory of Evidence of Dempster-Schafer in that the observers are not liars; from Fuzzy Logic in that no membership function is needed for each observation.
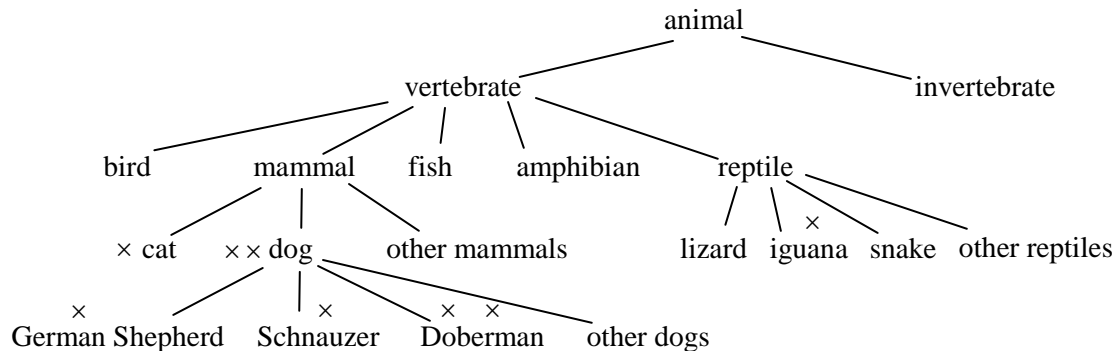
**Figure 1.** A hierarchy of symbolic values. It is a tree where every node is either a symbolic value or, if it is a set, then its descendants form a partition. Hierarchies make possible to compute the confusion conf($r, s$) that results when value $r$ is used instead of $s$, the true or intended value. The confusion is the number of *descending* links in the path from $r$ to $s$, divided by the height of the hierarchy. For instance, conf(dog, Doberman) = 1/4, conf(Doberman, dog) = 0, conf (Doberman, German Shepherd) = 1/4, conf (Doberman, iguana) = 2/4, conf(iguana, Doberman) = 3/4. conf $\in$ [0, 1]. Refer to Section 2. Values marked with $\times$ refer to Section 3

## 1. Previous work and problem statement

When several measurements are performed on the same variable (for instance, the length of a table), it is possible to obtain the most likely value ($\mu$=3.25m, the average length) as well as the dispersion of these measurements ($\sigma$, the variance), perhaps disregarding some outliers. For quantitative measurements we know how to take into account

contradicting facts, and we do not regard them necessarily as inconsistent. We just assume that the observers' gauges have different precisions or accuracies.

Let us now consider several observations on a non-numeric variable (such as pet_of_John_is) that ranges on qualitative values (such as dog, cat, German Shepherd, Schnauzer…) that can be arranged in a hierarchy (Figure 1). That is, observer 1 reports that John's pet is a dog, observer 2 reports that John's pet is a cat… Can we find the consensus value or most likely value for John's pet? The "centroid" or "average" of the reported pets? Can we find the "dispersion" (variance), discrepancy or degree of disagreement (inconsistency) of this bag[1] of values? Intuitively, this is the value that minimizes the sum of disagreements or discomforts for all the observers when they learn of the value chosen as the consensus value.

Section 3 of this paper solves the following

*Problem 1. Given a bag of observations reporting non-numeric values, how can me measure their inconsistency? What is the value that minimizes this inconsistency?* We shall call $r*$ this value and $\sigma$ the inconsistency that $r*$ produces.

Plausibility theory (Dempster-Schafer) solves this problem assuming that each observer has a given probability of lying, and that their observations are independent –they do not influence each other. We assume, instead, that all observers are truthful, so that the discrepancy in the values reported is due to the different methods used to perform the observation (observer 1 saw John's pet at a distance, observer 2 saw it at night, observer 3 examined its excrement…)

Logicians **[Hunter]** solve this problem by
(a) declaring that, since dog ≠ cat ≠Doberman ≠ …, the set is inconsistent, and the sentence (John's cat is a dog) ∧ (John's cat is a pet) ∧ …  evaluates to F; no agreement is possible; or
(b) postulating a (small) set of predicates that must all be true for this set of observations to be consistent, and declaring that the degree of inconsistency of the set is the percentage of predicates that become false; or
(c) using paraconsistent logic [   ]  falta aquí; or
(d) using non-monotonic logic [  ] falta aquí.
Fuzzy logic [  ] solves the problem assigning to each observer a fuzzy predicate or membership function, such as { (cat, ½), (dog, 1/6), (Doberman, 1/3) } (Figure 2) and then combining the function with the rules of fuzzy AND [dar un ejemplo]
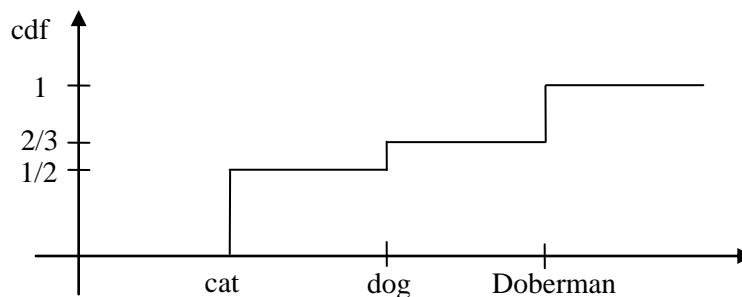


**Figure 2.** Membership function for observer 1, useful in the Fuzzy Logic approach

---

[1] A bag is a set where repeated elements are allowed.

Our solution uses hierarchies of qualitative values and the confusion conf($r$, $s$), explained in next section. From these, the inconsistency among a set of observations and their most likely value are defined (Section 3). Thus, Section 3 solves Problem 1. Section 4 solves Problem 1 when more than one answer are possible (John may have several pets!); that is, Section 4 solves

*Problem 2.* (informal statement) *Given a bag of non-numeric observations, what are the values that jointly minimize the inconsistency of the bag?*

The solution to this problem is a small set of "centers of gravity" $r^*_i$ or most likely pets, in our example.

It may turn out that the outliers in the bag of observations are clearly discernible and thus they may be safely discarded (as "observation errors"), so they do not contaminate neither $r^*$ nor $\sigma$. The following problem is meaningful:

*Problem 3. Given a bag of non-numeric observations (reported facts), what is the fact that reduces most the inconsistency of the remaining bag, if such fact is expunged?*

That is, find the most conspicuous outlier. Section 5 does that.
The problem solved in Section 6 is how to incrementally compute the inconsistency:

*Problem 4. Given $r^*$ and $\sigma$ for a bag of non-numeric observations, how do they change when a new observation is added?*

At times, observations come in bundles; they are agglutinations of simple observations, which can not be taken apart. For instance, witnesses of a crime give the following depositions:
Observation 1 = The murderer is a black, tall man with a tattoo in his body;
Observation 2 = The assassin is a Caucasoid, short man with a tattoo in his left finger;
Observation 3 = The killer is a Chinese, short woman with a tattoo in her right shoulder;
Observation 4 = The executor is a Japanese woman without a tattoo.
Notice that the facts of each observation can not be unglued or separated from each other. For instance, nobody observed a Caucasoid, tall woman with a tattoo in her right shoulder. The problem solved in Section 7 is

*Problem 5. Given a bag of objects described by a finite number of qualitative variables, find the object best considered as the "consensus" or "center of gravity" of such collection of objects.* An additional problem to be solved is

*Problem 5a. Which of the objects of the bag is most consistent with a given predicate? And with a given set of predicates?*

Finally, we consider the case when observations change with time. We now have a set of sequences of observations ordered on increasing dates (timestamps), such as
Sequence 1 = (John is in the street), (John is bleeding), (John is sick), (John goes into an ambulance), (John is in the hospital), (John is well), (John is at home)

Sequence 2 = (John is in the street), (John is walking), (John goes into a taxi), (John is well), (John is at home)
Sequence 3 = (John is in the street), (John is sick), (John is dying), (John is in a car), (John is at the cemetery).

We also have a set of dynamic models $M_1$, ..., $M_k$ (say, each is a finite state machine, or FSM) that describe different behaviors or processes. For instance, $M_1$ could be the FSM of Figure 3. Sequences may have different numbers of observations. The times $t_1 < t_2 \ldots < t_n$ in which observations of some sequence were taken, need not be identical to the times in which observations of another sequence were taken. For this reason, some state(s) of a given model may not have been observed by some observer –a missing observation. The problem to solve (in Section 8) is:

*Problem 6. Given a bag of sequences of observations and a set of dynamic models, which is the model that best fits (it has the lowest inconsistency) the bag?*

Section 9 provides discussions, conclusion and suggestions for further work.

**Figure 3.** Finite state machine $M_1$ gives one of the possible sequences of events that could occur to John (for Problem 6) falta la figura

## 2. Measuring the confusion among two qualitative values and among two objects

Here we extract from our work in [Levachkine & Guzman 2005 and 2007]. How close are two numeric values $v_1$ and $v_2$? The answer is $|v_2 - v_1|$. How close are two symbolic values such as cat and dog? The answer comes in a variety of similarity measures and distances, some of which are discussed in [Levachkine & Guzman 2007]. The hierarchies introduced in Figure 1 allows us to define the confusion conf($r, s$) on two symbolic values. We assume that the observers of a given fact (such as the pet of John) share a set of common vocabulary, best arranged in a hierarchy. This hierarchy can be regarded as the "common terminology"[2] of the observers, their *context.* Other observers may share a different context, that is, another hierarchy. The function conf will open the way to evaluate in Section 3 the inconsistency among a bag of symbolic observations.

What is the capital of Germany? *Berlin* is the correct answer; *Frankfurt* is a close miss, *Madrid* a fair error, and *sausage* a gross error. What is closer to a *cat*, a *dog* or an *orange*? Can we measure these errors and similarities? Can we retrieve objects in a database that are

---

[2] If the symbolic values become full *concepts,* it is best to use an *ontology* instead of a *hierarchy* to place them. [Cuevas & Guzman].

close to a desired item? Yes, by arranging these symbolic (that is, non-numeric) values in a hierarchy. More precisely, qualitative variables take symbolic values such as *cat, orange, California, Africa.* These values can be organized in a hierarchy *H*, a mathematical construct among these values. Over *H*, we can define the function *confusion* resulting when using a symbolic value instead of another.

*Definition.* For *r, s* ∈ H, the **absolute confusion** of using *r* instead of *s,* is
CONF(*r, r*) = CONF(*r,* any ascendant of *r*) = 0;
CONF(*r, s*) = 1 + CONF(*r,* father_of(*s*)).

 To measure CONF, count the descending links from *r* (the replacing value) to *s* (the intended or real value). CONF is not a distance, nor an ultradistance.

 We can normalize CONF by dividing into *h,* the height of H (the number of links from the root of H to the farthest element of H), yielding the following

*Definition.* The **confusion** of using *r* instead of *s* is
conf(*r, s*) = CONF(*r, s*)/h.

 Notice that $0 \leq conf(r, s) \leq 1$. It is not symmetric: conf(*r, s*) ≠ conf(*s, r*), in general.
 Example. In the hierarchy of Figure 1, conf(cat, mammal)=0; conf(cat, dog)=1.

## 2.1 Confusion among two objects

 For us, an object is defined by a list of qualitative values. It is also possible to handle a mixture of qualitative and numeric values. Thus, O = (tall, Mexico, iguana), meaning perhaps that object O is tall, lives in Mexico and has an iguana as pet. If O' = (tall, American Continent, reptile), then we can measure the confusion of using O' instead of O, by just adding [Levachkine & Guzman 2007] the confusions that their respective properties provoke, thus: conf(O, O') = conf(tall, tall) +conf(Mexico, American Continent) +conf(iguana, reptile) = 0 + 0 + 0 = 0, whereas conf(O', O) = conf(tall,tall) + conf(American Continent, Mexico) + conf(reptile, iguana) = 0 + 2/3 + 1/4 = 0.91, using suitable hierarchies such as those of Figure 1 and Figure 4. These hierarchies represent the *context* or common vocabularies of O, O' and other objects; without them their closeness can not be ascertained.
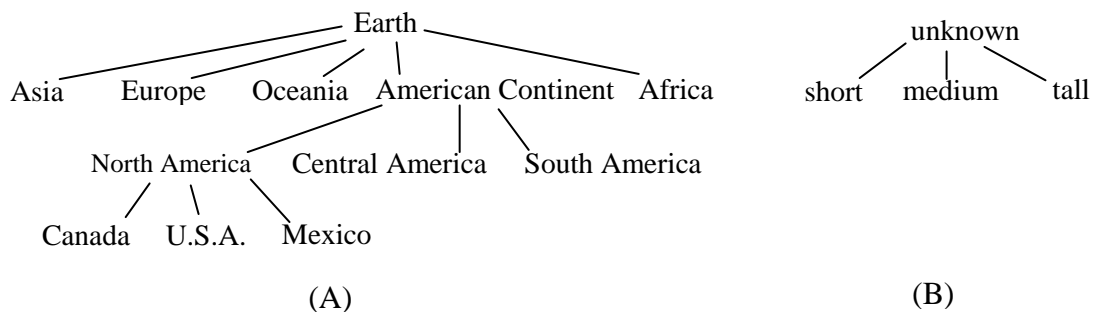


**Figure 4.** Hierarchies for places to live (A) and height of persons (B)

 To obtain a normalized confusion among two objects, we divide by their number of properties. Thus, we extend conf to work on objects described by *m* properties, as follows:

$$\text{conf(O, O')} = (1/m) \sum_{i=1}^{m} \text{conf}(o_i, o_i') \qquad \text{where } o_i \text{ is the } i\text{th property of O, and so for } o_i'.$$

This formula will be useful when computing in Section 7 the inconsistency produced by a bag of objects.

## 3. Measuring the degree of inconsistency

*Problem 1. Given a bag of observations reporting non-numeric values, how can me measure their inconsistency? What is the value that minimizes this inconsistency?*

Assume we have a bag of observations about the pet of John: observers report that John has a dog, a dog, a cat, a German Shepherd, a Schnauzer, a Doberman, a Doberman, an iguana, respectively. These observations are marked with × in Figure 1. We wish to determine how divergent or discrepant they are –how inconsistent they are. Notice that inconsistency is a property of a bag of observations.

We first measure the total confusion that occurs when a given value $r$ (one of the ×) is used instead of all the reported values (×). For instance, for $r$ = cat, we obtain

Total confusion when cat is selected as the "representative" of the bag of observations = conf(cat,dog) + conf(cat,dog) + conf(cat,cat) + conf(cat,German Shepherd) + conf(cat,Schnauzer) + conf(cat,Doberman) + conf(cat,Doberman) + conf(cat,iguana) = ¼+¼+0+½+½+½+½+½=3. The first term of the sum, conf(cat,dog)= ¼, means that the observer that reported "dog" will be slightly annoyed (conf=¼) when he finds that the representative is "cat" Similarly, conf(cat, German Shepherd) = ½ means that the observer that saw a German Shepherd will be at discomfort = ½ when he finds that the consensus is cat. The sum of confusions thus measures the total disagreement with the chosen representative value.

Total confusion if dog were the "representative" of the observations = conf(dog,dog) +conf(dog,dog) +conf(dog,cat) +conf(dog,German Shepherd) +conf(dog,Schnauzer) +conf(dog,Doberman)+conf(dog,Doberman)+conf(dog,iguana)=0+0+¼+¼+¼+¼+¼+½=7/4.

Total confusion for German Shepherd as the representative = 0+0+¼+0+¼+¼+¼+½ = 3/2.

Total confusion for Schnauzer as the representative = 0+0+¼+¼+0+¼+¼+½ = 3/2.

Total confusion for Doberman as the representative = 0+0+¼+¼+¼+0+0+½ = 5/4.

Total confusion for iguana as the representative = ½+½+½+3/4+3/4+3/4+3/4+0 = 9/2.

It makes sense to take as the best representative (consensus value) the animal that minimizes the total confusion. Such animal is Doberman. This is the best "consensus", because it minimizes the confusion or "discomfort" of the observers when they see that Doberman was selected, instead of their observed pet. We can call this the "centroid" of the observed facts (× in Figure 1), the $r^*$ of Problem 1. Also,

Average confusion = total confusion / number of facts = (5/4)/8 = 5/32. This average confusion is called the *inconsistency* $\sigma$ of the bag of observations. It is the average confusion produced by its centroid $r^*$. Therefore, we have the following

*Definitions.* The *centroid* or *consensus* $r^*$ of a bag B of observations reporting qualitative values $\{s1, s2, \ldots, sk\}$ is the $r_j \in$ B that minimizes

$$\sum_{i=1}^{k} \mathrm{conf}(r_j, s_i) \qquad \text{for } j = 1,\ldots, k$$

The *inconsistency* of the bag is the minimum that such $r^*$ produces, divided by k:

$$\sigma = (1/k) \min_{j \in [1,k]} \sum_{i=1}^{k} \mathrm{conf}(r_j, s_i) = (1/k) \sum_{i=1}^{k} \mathrm{conf}(r^*, s_i)$$

Remarks. (A) The above centroid and inconsistency are the solutions to Problem 1. The inconsistency $\sigma \in [0, 1)$. (B) The consensus $r^* \in \{s1, s2, \ldots, sk\}$. (C) There may be more than one value $r^*$ that minimizes the total confusion. (D) To compute the inconsistency of a bag, we resort to finding $r^*$. (E) $r^*$ is not necessarily the most popular value (the mode). (F) The *least common ancestor* (vertebrate in our example) produces a total confusion larger or at best equal than the total confusion produced by $r^*$. (G) If we could discard *outliers* (Section 5 shows how)*, then iguana could be discarded; that will still produce the consensus $r^*$ = Doberman, but now with a $\sigma$ = 1/14, a much tighter result. [In general, the new consensus may shift].

Examples. For bag1 = {animal, vertebrate, bird, mammal, cat, dog, dog, iguana, German Shepherd} (marked with × in Figure 5), the centroid $r^*$ is German Shepherd, and the inconsistency of the bag is (4/9)/5 = 4/45. For bag2 = {animal, amphibian, amphibian, reptile, reptile, snake} marked with • , $r^*$ = snake, $\sigma$ = (2/6)/5 =1/15. Taking into account all the observations × and •, we obtain for bag1∪bag2 a consensus $r^*$ = German Shepherd with $\sigma$ = (10/15)/5 = 2/15.



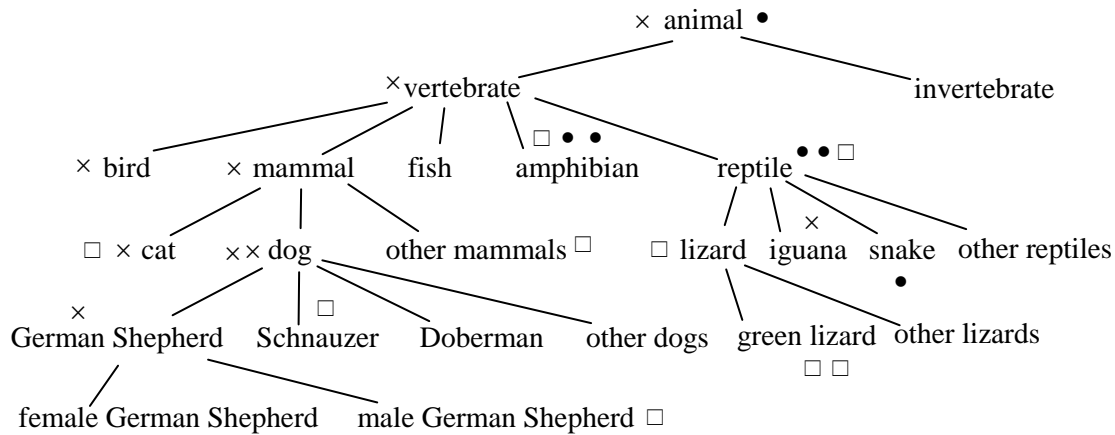**Figure 5.** The consensus of the observations marked (×) is $r^*$ = German Shepherd, with inconsistency $\sigma$ = 4/45. For observations marked (•), $r^*$ = snake, $\sigma$ = 1/15. For observations marked (□), $r^*$ = green lizard, $\sigma$ = (12/9)/5 = 4/15

Notice that we have found a way of adding (and averaging) apples and oranges, and a way ($\sigma$) to find out how disperse or divergent a bag of symbolic values is.

Properties of $r^*$ (consensus) and $\sigma$ (inconsistency). Stated without proof.

- $r*$ and σ depend on the context of use –represented by the hierarchy employed.
- Could the consensus be outside B? Given a bag B, there is no other value $r+ \neq r*$ in the hierarchy that provokes a lower value for σ. For instance, for bag1 in Figure 5, all ascendants of German Shepherd will yield inconsistencies larger than 4/45. We could have used Female German Shepherd instead of German Shepherd, which also yields σ=4/45. But it will be strange to report that the consensus is Female German Shepherd, since no observer reported this value! Hence, in the definition we force $r* \in$ B.
- The more specialized, the better. Consensus tends to go to the precise values (those deep in the hierarchy), unless of course overruled by several other less-precise values. For instance, the consensus of $\{s_1, s_2,..., s_k\}$ where $s_i < s_j$ ($s_i$ is a descendant of $s_j$) whenever i < j, is $s_1$. Example: In figure 4, the consensus of {Doberman, dog, mammal} is Doberman.

## 4. When several consensus values are possible

A person is born at just one place, and has precisely one mother, whereas he can have several friends, or several pets. Some properties, then, can have more than one value. In such case, bag B can have several centroids, and Problem 1 could be stated as

*Problem 2.* (tentative) *Find k different consensus $r_1*, r_2*, ...,r_k*$ such that the inconsistency defined as*

$$\sigma = (1/|B|) \sum_{i=1}^{|B|} conf(s_i, r_j*) \text{ is smallest, where the confusion of each } s_i \in B \text{ is measured}$$

*against the consensus $r_j*$ that has smallest* $conf(s_i, r_j*)$.

Problem 2 has a trivial solution: take each different member of B as one of the consensus $r_j*$. Then, the number of consensus $k$ is ≤ |B|. It is easy to see that $\sum_{i=1}^{|B|} conf(s_i, r_j*)$ is 0, since the confusion of any $s_i$ with itself is 0.

The number of consensus of the trivial solution presented above could be further reduced if we delete from them each consensus $r_j*$ that is ascendant of some other consensus $r_m*$. That is because the observation $r_j*$ could still contribute with confusion 0 since it will select now as "its" consensus the value $r_m*$, thus: $conf(r_m*, r_j*) = 0$.

We must penalize solutions with many centroids. One way is to introduce a penalty $0 \leq p \leq 1$ so that the function to minimize in Problem 2 is a combination of the number of centroids and the inconsistency.

*Problem 2. For a given $p \in [0,1]$, find the k different consensus $r_1*, ..., r_k*$ that minimize*

$$pk + (1-p)(1/|B|) \sum_{i=1}^{|B|} conf(s_i, r_j*) \text{ where the confusion of each } s_i \in B \text{ is measured against the}$$

*consensus $r_j*$ that has smallest* conf($s_i$, $r_j*$). A large $p$ will render few clusters, a small $p$ will produce smaller inconsistencies.

It may be advisable to get rid of the subjectivity of selecting a suitable $p$. One way could be to watch how the inconsistency drops as $k$ grows, and to select the $k$ that causes the first considerable drop, if such $k$ exists. For instance, in Figure 5, a sizable drop in the inconsistency of bag1∪bag2 exists going from $k=2$ to $k=3$ (see Table 1), so we select $k=3$.

**Table 1.** Centroids and inconsistency of bag1∪bag2, which consists of all the observations marked × and • in Figure 5. There is a big drop in inconsistency from 0.12 for two centroids to 0.09 for three centroids, so we select the three centroids German Shepherd, amphibian, iguana as a good "consensus"

| Number of centroids | Centroids  checar si son los buenos | inconsistency |
|---|---|---|
| 1 | German Shepherd | $2/15 = 0.13$ |
| 2 | German Shepherd, snake | $(4/9 + 1/6)/5 = 0.12$ |
| 3 | German Shepherd, amphibian, iguana | $(2/7 + 0/2 + 1/6)/5 = 0.09$ |
| 4 | German Shepherd, cat, amphibian, snake | $(0/3 + 1/6 + 0/2 + ¼)/5 = 0.08$ |
| 5 | German Shepherd, cat, bird, amphibian, iguana | $(0/7 + 0/1 + 0/1 + 0/2 + ¼)/5 = 0.05$ |
| 6 | German Shepherd, cat, bird, amphibian, iguana, snake | $0/7 + 0/1 + 0/1 + 0/2 + 0/3 + 0/1 = 0$ |

What we are doing is clustering the qualitative values of the bag, and finding the consensus or centroid inside each cluster.

Remark. Once we know how many clusters or centroids we wish to have (this number is $k$), it is a well defined minimization problem to find the $k$ centroids of a bag B of qualitative values: we want to minimize the resulting inconsistency for $k=3$. In the worst case, we could evaluate all the partitions of B into three sub-bags. There exist methods and heuristics to find reasonable clusters without exhaustive search, but we shall not delve into them. See [Adriana Jiménez] for one of these.

Remark. When $p = 0$, no importance is given to the number of clusters, thus the inconsistency of any bag is 0 and the centroids are found by the trivial solution to Problem 2. When $p = 1$, no importance is given to the value of the inconsistency; one centroid is found through the solution to Problem 1.
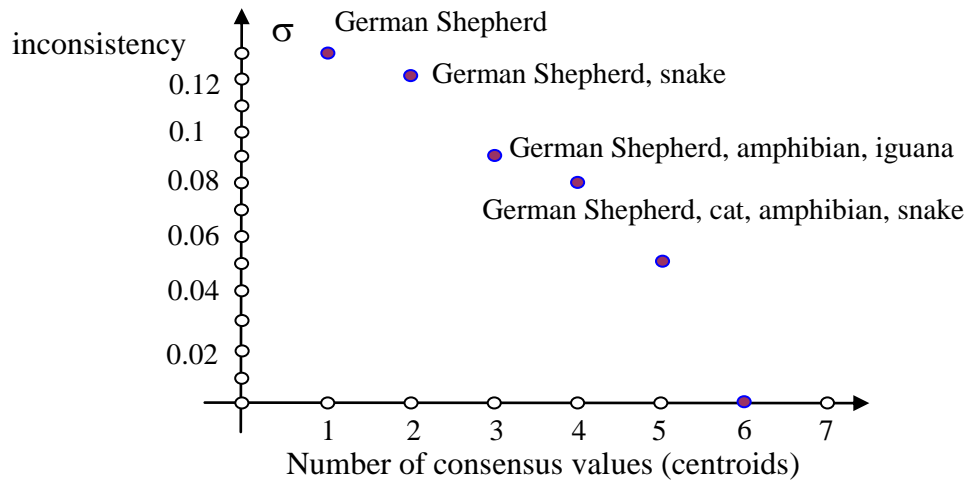
**Figure 6.** The inconsistency (Table 1) decreases as more consensus values are allowed as representatives of bag1∪bag2 (observations marked × and • in Figure 5)

# 5. Detecting the most prominent outlier

Here we solve

*Problem 3. Given a bag of non-numeric observations, what is the observation that reduces most the inconsistency when removed from the bag?*

The solution is easy: remove each different value from the bag, compute $(r^*, \sigma)$ for the resulting bag, take the value that minimizes $\sigma$. We must compute only one centroid (not several, as in Problem 2) for the resulting bag.

*Solution.* The most conspicuous outlier $r_O$ = the value $r_i$ that minimizes the inconsistency of $B - \{r_i\}$, where $r_i$ is each different value of B. Report that the inconsistency has dropped from $\sigma$ to $\sigma_O$. *Caution:* If there are several repeated values inside the bag, removing one of them will not remove the others. *Hint:* try first candidates away from the consensus. *Remark.* There may be more than one $r_O$. *Remark:* removal of $r_O$, in addition to producing a smaller redundancy, may result in a new centroid (for the smaller bag).

Examples. ♦ The outermost outlier for the bag of objects marked with × in Figure 1 is iguana; expunging it produces Doberman as consensus, with a new inconsistency of $(2/7)/4 = 1/14$, as opposed to 5/32 when undeleted (as computed in Section 3). ♦ For bag1 (values marked × in Figure 5), the centroid is German Shepherd, the inconsistency is 4/45 (according to Figure 5). Its most conspicuous outlier is iguana; removing this value reduces the inconsistency of the smaller bag to $(2/8)/5=1/20$; the new centroid is still German Shepherd. ♦ The bag2 of Figure 5 (values marked with •) has an inconsistency of 1/15, with centroid = snake. Its most prominent outlier is amphibian. Without it, the smaller bag2 drops its inconsistency to $(1/5)/5 = 1/25$, the new centroid is still snake.

# 6. How to incrementally compute the inconsistency

Once we have found the inconsistency of a bag, how will it be altered when a new observation is added? This is

*Problem 4. Given r\* and σ for a bag of non-numeric observations, how do they change when a new observation s is added?*

It is not possible, just knowing *s*, *r\** and σ for a bag, to compute its new σ when a new observation *s* is added. We must also know |B|, the number of observations in the bag. Then, the solution can be found by computing the new total confusion σ |B| + conf(*r\**, *s*) and dividing it among the new number of elements of B, |B| +1. That is:

$$\sigma_{NEW} = [\sigma\ |B| + conf(r^*,\ s)] / [|B| +1]$$

(APPROXIMATION 1)

$$r^*_{NEW} = r^*$$

Unfortunately, this new inconsistency is computed against the old consensus *r\**, which may not be the new consensus anymore. Approximation 1 is suitable when the number |B| of observations is large (say, >10), and when the new observation *r* it not too far from *r\** (say, conf(*r\**, *r*) < σ). In these cases, it is reasonable to assume that the consensus will remain unchanged. Exact computation of $r^*_{NEW}$, $\sigma_{NEW}$ requires to know all the observations of bag B, and then apply to B ∪ {*s*} the solution given in Section 3 to Problem 1.

Examples. Table 2 gives the new inconsistency and consensus when a new value *bird* is added to some bag.

**Table 2.** Incremental change of the inconsistency and the consensus when a new value *s* = bird is added

|  | Old consensus and inconsistency | New inconsistency (by approximation 1) | New consensus and inconsistency (correct values) | Was Approx. 1 good? |
|---|---|---|---|---|
| Bag × of Figure 1 | Doberman, 5/32 | (5/4 + ¼)/9=3/18 | Doberman, (6/9)/4 =3/18 | ✓ |
| Bag × of Figure 5 | German Shepherd, 4/45 | (4/5 + 1/5)/10 = 1/10 | German Shepherd, (5/5)/10= 1/10 | ✓ |
| Bag ● of Figure 5 | snake, 1/15 | (6/15+1/5)/7=3/35 | snake, (3/7)/5 = 3/35 | ✓ |
| Bag {×}∪{●} of Figure 5 | German Shepherd, 2/15 | (2+1/5)/16 = 0.137 | German Shepherd, (11/16)/5 = 0.137 | ✓ |
| Bag {□} of Figure 5 | green lizard, 4/15 | (9*4/15 + 1/5)/10 = 0.26 | green lizard, (13/5)/10 =0.26 | ✓ |

# 7. Consensus and inconsistency among a bag of objects

The problems to be solved in this section are:

*Problem 5. Given a bag of objects described by a finite number of qualitative variables, find the object best considered as the "consensus" or "center of gravity" of such collection of objects.* An additional problem to be solved is

Often, observers report several properties of a given object, so that the observations come in bundles. For instance, each observer in Table 3 observed the same object and wrote down in a row of the table its properties (left half of the table).

With the help of the confusion conf(O, O') that results when object O is used instead of object O' (Section 2.1), and with the help of suitable hierarchies that define the context of the qualitative values of the objects (we use for this example the hierarchies of Figure 1 and Figure 4), we can obtain the confusion between each pair of objects. For instance, conf($O_2$, $O_4$) = 1/3 *[conf(tall, medium) +conf(American Continent, American Continent + conf(reptile, vertebrate)] = 1/3 * [1+0+0] = 0.33, whereas conf($O_4$, $O_2$) = 1/3*[1+0+1/4] = 0.42. In this way, the right side of Table 3 is obtained.

**Table 3.** The object reported by each observer appears as one row in the left side of this table. The right side represents the confusion when object *r* (a row) is used instead of object *c* (a column); for instance, conf($O_1$, $O_3$) = [conf(tall, short) + conf(Mexico, Canada) + conf(iguana, Schnauzer)]/3 = [1+1/3+3/4]/3 = 0.69

| Observer | height | Place of living | Pet he has | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | tall | Mexico | iguana | 0 | 0 | 0.69 | 0.33 | 0.11 | 0 |
| 2 | tall | American Continent | reptile | 0.31 | 0 | 0.81 | 0.33 | 0.53 | 0 |
| 3 | short | Canada | Schnauzer | 0.61 | 0.42 | 0 | 0.33 | 0.11 | 0 |
| 4 | medium | American Continent | vertebrate | 0.72 | 0.42 | 0.81 | 0 | 0.19 | 0 |
| 5 | medium | Africa | mammal | 0.83 | 0.53 | 0.83 | 0.11 | 0 | 0 |
| 6 | unknown | Earth | animal | 0.92 | 0.61 | 1 | 0.53 | 0.61 | 0 |

Armed with the confusion among two objects, we can now find the inconsistency and the centroid of a bag of *objects* using the same formulas given Section 3 for values, being careful to use in them conf for objects.

Example: for bag {1,2,2,3} of objects in Table 3, the consensus is object 1 with an inconsistency of [conf(1,1)+conf(1,2)+conf(1,2)+conf(1,3)]/4 = (0+0+0+0.69)/4 = 0.17. [3]

Example: for bag {2, 3, 4, 4, 5}, the consensus is object 3 with an inconsistency of (0.42 + 0 +0.33 +0.33 + 0.11)/5 = 1.19/5 = 0.238.

Example: for bag {2, 4, 6}, the consensus is object 2 with an inconsistency of (0 + 0.33 + 0)/3 = 0.11.

The following formulas formalize the results.

The *centroid* or *consensus O\** of a bag B of objects {O1, O2, …, Ok}described by qualitative values, is the object $O_j \in B$ that minimizes

$$\sum_{i=1}^{k} conf(O_j, O_i) \quad \text{for } j = 1,..., k$$

---

[3] If we had chosen object 2 as the centroid of the bag, its inconsistency would be (0.31 + 0 + 0 + 0.81)/4 = 0.28. If we had chosen object 3 as the centroid of the bag, its inconsistency would be (0.61 + 0.42 + 0.42 + 0)/4 = 0.36. Thus, object 1 with inconsistency 0.17 has the lowest inconsistency; therefore, object 1 is the centroid or consensus of the values of the bag {1, 2, 2, 3}.

The *inconsistency* of the bag is the minimum that such *O\** produces, divided by k:

$$\sigma = (1/k) \min_{j \in [1,k]} \sum_{i=1}^{k} \mathrm{conf}(O_j, O_i) = (1/k) \sum_{i=1}^{k} \mathrm{conf}(O^*, O_i)$$

The objects in the bag are all described by the same properties or attributes (such as place of origin, color of hair, religion…); the *values* of such properties will vary, of course, from object to object.

It is possible to solve Problem 5a: to find which of the objects of a set is most consistent with a given predicate P. That is, P is a predicate that evaluates to a number between 0 (false) and (1) true; when applied to an object *O*, P(*O*) evaluates the *inconsistency* between *O* and the predicate P. This problem is solved in [Levachkine & Guzman 2007], where it is called "object *O* fulfils predicate P with confusion ε."

# 8. Inconsistency for time-varying values – fitting a dynamic model

falta

# 9. Discussion and conclusion

Falta

**References**

Fuzzy logic [  ]
[Adriana Jiménez]  ph d thesis
[Hunter]
[Levachkine & Guzman 2005] Sergei Levachkine, Adolfo Guzman-Arenas, Victor Polo de Gyves (2005) The semantics of confusion in hierarchies: from theory to practice. In *Contributions to ICCS 05 13th International Conference on Conceptual Structures: common semantics for sharing knowledge*, July 18-22 2005, Kassel, Germany. 94-107
[Levachkine & Guzman 2007] Serguei Levachkine, A. Guzman-Arenas (2007) Hierarchy as a new data type for qualitative variables. *Journal Expert Systems with Applications* **32**, 3, June 2007
non-monotonic logic [  ]
paraconsistent logic [   ]