# Ugly duckling theorem

The **Ugly Duckling theorem** is an argument asserting that classification is impossible without some sort of bias. It is named for Hans Christian Andersen's famous story of "The Ugly Duckling." It gets its name because it shows that, all things being equal, an ugly duckling is just as similar to a swan as two swans are to each other, although it is only a theorem in a very informal sense. It was proposed by Satosi Watanabe in 1969.[1]

## Basic idea

Suppose there are n  things in the universe, and one wants to put them into classes or categories. One has no preconceived ideas or biases about what sorts of categories are "natural" or "normal" and what are not. So one has to consider all the possible classes that could be, all the possible ways of making sets out of the n  objects. There are $2^n$ such ways, the size of the power set of n  objects. One can use that to measure the similarity between two objects: and one would see how many sets they have in common. However one can not. Any two objects have exactly the same number of classes in common if they are only distinguished by their names with one another, namely $2^{n-1}$ (half the total number of classes there are). To see this is so, one may imagine each class is a represented by an n-bit string (or binary encoded integer), with a zero for each element not in the class and a one for each element in the class. As one finds, there are $2^n$ such strings.

As all possible choices of zeros and ones are there, any two bit-positions will agree exactly half the time. One may pick two elements and reorder the bits so they are the first two, and imagine the numbers sorted lexicographically. The first $2^n/2$ numbers will have bit #1 set to zero, and the second $2^n/2$ will have it set to one. Within each of those blocks, the top $2^n/4$ will have bit #2 set to zero and the other $2^n/4$ will have it as one, so they agree on two blocks of $2^n/4$ or on half of all the cases. No matter which two elements one picks. So if we have no preconceived bias about which categories are better, everything is then equally similar (or equally dissimilar). The number of predicates simultaneously satisfied by two non-identical elements is constant over all such pairs and is the same as the number of those satisfied by one. Thus, some kind of inductive bias is needed to make judgements; i.e. to prefer certain categories over others.
(A possible way to proceed is however correspondence analysis).

## As a statement about Boolean functions

Let $x_1, x_2, \ldots, x_n$ be a set of vectors of $k$ booleans each. The ugly duckling is the vector which is least like the others. Given the booleans, this can be computed using Hamming distance.

However, the choice of boolean features to consider could have been somewhat arbitrary. Perhaps there were features derivable from the original features that were important for identifying the ugly duckling. The set of booleans in the vector can be extended with new features computed as boolean functions of the $k$ original features. The only canonical way to do this is to extend it with *all* possible Boolean functions. The resulting completed vectors have $2^k$ features. The Ugly Duckling Theorem states that there is no ugly duckling because any two completed vectors will either be equal or differ in exactly half of the features.

Proof. Let x and y be two vectors. If they are the same, then their completed vectors must also be the same because any Boolean function of x will agree with the same Boolean function of y. If x and y are different, then there exists a coordinate $i$ where the $i$ -th coordinate of $x$ differs from the $i$ -th coordinate of $y$. Now the completed features contain every Boolean function on $k$ Boolean variables, with each one exactly once. Viewing these Boolean functions as polynomials in $k$ variables over GF(2), segregate the functions into pairs $(f, g)$ where $f$ contains the $i$ -th coordinate as a linear term and $g$ is $f$ without that linear term. Now, for every such pair $(f, g)$, $x$ and $y$ will agree on exactly one of the two functions. If they agree on one, they must disagree on the other and vice versa. (This proof is believed to be due to Watanabe.)

# Notes

[1] Watanabe, Satosi (1969). *Knowing and Guessing: A Quantitative Study of Inference and Information*. New York: Wiley. pp. 376–377.

# Article Sources and Contributors

**Ugly duckling theorem**  *Source*: http://en.wikipedia.org/w/index.php?oldid=455288346  *Contributors*: Arthur Rubin, Avalon, CBM, Clsn, FMB, Giftlite, Gregbard, Incnis Mrsi, JRSpriggs, Jensbn, JonHarder, Lambiam, Michael Hardy, Pallab1234, RDBury, Rjwilmsi, Silly rabbit, Tassedethe, ThomHImself, Took, U664003803, Zvika, 12 anonymous edits

# License