

Word Representation: Explanation of an Empirical Formula

Leonel Escapita, Diana Licon, Madison Anderson,
Diego Pedraza, and Vladik Kreinovich
Department of Computer Science, University of Texas at El Paso
El Paso, Texas 79968, USA
mranderson2@miners.utep.edu, laescapita@miners.utep.edu
dlicon2@miners.utep.edu, dapedraza@miners.utep.edu, vladik@utep.edu

Formulation of the Problem. Computer representation of natural-language words should reflect how related these words are. This relation can be described, e.g., by the number X_{ij} of times when word i appears in the context of word j . It is desirable to find out how such characteristics depend on the properties of the words, this way we will be able to predict, e.g., X_{ij} for pairs (i, j) for which we do not know X_{ij} . At present, these characteristics are determined by training a neural network.

There is also a reasonably good approximate analytical formula $\ln(X_{ij}) \approx b_i + \tilde{b}_j + w_i \cdot \tilde{w}_j$. Here, b_i and \tilde{b}_j are numbers, w_i and \tilde{w}_j are vectors, and $a \cdot b$ is dot (scalar) product. The values b_i , \tilde{b}_j , w_i and \tilde{w}_j can be found by using the Least Squares method:

$$J \stackrel{\text{def}}{=} \sum_{i,j} f(X_{ij}) \cdot (b_i + \tilde{b}_j + w_i \cdot \tilde{w}_j - \ln(X_{ij}))^2 \rightarrow \min.$$

The efficiency of this method depends on the appropriate choice of the weight function $f(X)$. Empirical data shows that the most efficient is power law $f(X) = X^a$. How can we explain this empirical fact?

Our Explanation. The values X_{ij} depend on the size of the corpus. In general, if we consider a λ times larger corpus, we will get new values which are close to $\lambda \cdot X_{ij}$. The word representation should depend only on the words, not on corpus size. So, the resulting representation should not change if we replace X_{ij} with $\lambda \cdot X_{ij}$. Of course, if we replace X_{ij} with $\lambda \cdot X_{ij}$, the weights will change. However, this does not necessarily mean that the resulting representations will change.

Namely, for any $c > 0$, optimizing any function J is equivalent to optimizing the function $c \cdot J$. Example: the richest person on Earth is the richest whether we count his richness in dollars or in pesos. If we replace $f(X)$ with $c \cdot f(X)$, we will get $c \cdot J$ instead of J , so we will get the same representations w_i . Thus, we can have $f(\lambda \cdot X) = c \cdot f(X)$, and the resulting representation will be the same. So, the invariance with respect to corpus size means that for every $\lambda > 0$, there exists c depending on λ for which $f(\lambda \cdot X) = c(\lambda) \cdot f(X)$. It is known that all measurable solutions to the functional equation $f(\lambda \cdot X) = c(\lambda) \cdot f(X)$ are power laws $f(X) = A \cdot X^a$. As we have mentioned, the use of such weights is equivalent to using the weights $f(X) = X^a$.

This explains why the power law weights work the best: power law weights are the only ones for which the resulting representation does not depend on the corpus size.

How to Prove the Result about the Functional Equation. This result is easy to prove when the function $f(X)$ is differentiable. Suppose that $f(\lambda \cdot X) = c(\lambda) \cdot f(X)$. If we differentiate both sides with respect to λ , we get $X \cdot f'(\lambda \cdot X) = c'(\lambda) \cdot f(X)$. In particular, for $\lambda = 1$, we get $X \cdot f'(X) = a \cdot f(X)$, where $a \stackrel{\text{def}}{=} c'(1)$, so $X \cdot \frac{df}{dX} = a \cdot f$. We can separate the variables if we multiply both sides by $\frac{dX}{X \cdot f}$, then we get $\frac{df}{f} = a \cdot \frac{dX}{X}$. Integrating both sides of this equality, we get $\ln(f) = a \cdot \ln(X) + C$. By applying $\exp(x)$ to both sides, we get $f(X) = \exp(a \cdot \ln(X) + C) = A \cdot X^a$, where $A \stackrel{\text{def}}{=} e^C$.