

Attention in machine learning: how to explain the empirical formula

Sobita Alam, Arman Hossain, Samin Islam,
Arin Rahman, and Vladik Kreinovich

Department of Computer Science, University of Texas at El Paso
500 W. University, El Paso, TX 79968, USA

salam@miners.utep.edu, arahman6@miners.utep.edu,
sislam3@miners.utep.edu, ahossain4@miners.utep.edu, vladik@utep.edu

1. Attention as a way to better classification

- In many practical situations, we have several objects.
- Each object i is characterized by vector $x_i = (x_{i,1}, \dots, x_{i,N})$ consisting of this object's numerical characteristics.
- For example, we have many picture of pets, and you want to classify them into cats and dogs.
- One of the difficulties is that objects within each class are different.
- For example, dogs can be large and small, of different breeds, etc.
- To make classification task easier, it is desirable:
 - to replace each specific vector x_i
 - with a weighted average $y_i = \sum_j w_{ij} \cdot x_j$ of all the objects x_j which are similar to x_i .
- This way, the role of individual characteristics will diminish, and the classification task will become easier.

2. Attention as a way to better classification (cont-d)

- A natural way to describe the closeness between the objects x_i and x_j is to use the usual metric $d(a, b) = \sqrt{\sum_k (a_k - b_k)^2}$.
- The smaller this distance, the larger should be the weight.
- So we must have $w_{ij} \sim f(d(x_i, x_j))$ for some decreasing function $f(v)$.
- The sum of the weights should be equal to 1, so we must have

$$w_{ij} = \frac{f(d(x_i, x_j))}{\sum_{\ell} f(d(x_i, x_{\ell}))}.$$

- This expression can be simplified if we take into account that overall, the values x_{ij} are reasonably random.
- In this case the value $x_i^2 = \sum_k x_{i,k}^2$ is close to some constant C (which is N time average of $x_{i,j}^2$).
- Then, $d^2(x_i, x_j) = x_i^2 + x_j^2 - 2x_i \cdot x_j \approx 2C - 2x_i \cdot x_j$.
- So, a decreasing function of $d(x_i, x_j)$ can be described as an increasing function of the dot product $x_i \cdot x_j$.

3. Attention as a way to better classification (cont-d)

- Thus,

$$w_{ij} = \frac{F(x_i \cdot x_j)}{\sum_{\ell} F(x_i \cdot x_{\ell})}.$$

- Empirical evidence shows that out of all increasing functions $F(v)$, functions $F(v) = \exp(\alpha \cdot v)$ work the best.
- How can we explain this empirical fact?

4. Our explanation

- It is based on the fact that measurements are noisy.
- So, a natural requirement is that the resulting values y_i should be affected by the noise as little as possible.
- What if we replace the original values $x_{i,j}$ with noisy values $\tilde{x}_{i,k} = x_{i,k} + n_{i,k}$ for some noise $n_{i,k}$ with 0 mean.
- Then the dot product $\tilde{x}_i \cdot \tilde{x}_j$ becomes $x_i \cdot x_j + x_i \cdot n_j + n_i \cdot x_j + n_i \cdot n_j$.
- The expected value of terms $x_i \cdot n_j$ is 0, so the only non-zero addition to the dot product is $E[n_i \cdot n_j]$.
- For local noise, this expected value is 0; however:
 - if the noise had a global component with mean square value m ,
 - then, on average, all dot products are increased by the same constant m .
- So, we want to find the function $F(v)$ for which adding a constant m to all dot product would not change the weights.

5. Our explanation (cont-d)

- In particular, for two objects, this means that

$$\frac{F(a+m)}{F(a+m)+F(b+m)} = \frac{F(a)}{F(a)+F(b)} \text{ for all } a, b, \text{ and } m.$$

- If we apply $1/z$ to both sides of this equality and subtract 1 from both sides, we get

$$\frac{F(b+m)}{F(a+m)} = \frac{F(b)}{F(a)}.$$

- Multiplying both sides by $F(a+m)/F(b)$, we get

$$\frac{F(b+m)}{F(b)} = \frac{F(a+m)}{F(a)} \text{ for all } a \text{ and } b.$$

- So, the ratio $F(a+m)/F(a)$ does not depend on a , it only depends on m : $F(a+m)/F(a) = g(m)$ for some function $g(m)$.
- Thus, $F(a+m) = g(m) \cdot F(a)$.

6. Our explanation (cont-d)

- It is known that the only increasing solution to this functional equation is $F(a) = c \cdot \exp(\alpha \cdot a)$.
- From the viewpoint of the weights $w_{i,j}$, it is equivalent to

$$F(a) = \exp(\alpha \cdot a).$$

- This is exactly what we needed to explain.

7. Comment: how to solve the functional equation

- To solve the functional equation, differentiate both sides by m and take $m = 0$.
- Then $F'(a) = g'(0) \cdot F(a)$, with $\alpha \stackrel{\text{def}}{=} g'(0)$, i.e.,

$$\frac{dF}{da} = \alpha \cdot F \text{ and } \frac{dF}{F} = \alpha \cdot da.$$

- Integrating, we get $\ln(F) = \alpha \cdot a + \text{const}$, so $F(a) = \text{const} \cdot \exp(\alpha \cdot a)$.

8. Acknowledgments

This work was supported in part:

- by the US National Science Foundation grants:
 - 1623190 (A Model of Change for Preparing a New Generation for Professional Practice in Computer Science),
 - HRD-1834620 and HRD-2034030 (CAHSI Includes),
 - EAR-2225395 (Center for Collective Impact in Earthquake Science C-CIES),
- by the AT&T Fellowship in Information Technology, and
- by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) Focus Program SPP 100+ 2388, Grant Nr. 501624329,