

Non negative Matrix factorization

We have a data matrix X containing non-negative entries. Also X has n rows and m columns. We can decompose X into matrices W and H whose shapes are $n \times k$ and $k \times m$. Thus, $X_{ij} \approx \sum_k W_{ik}H_{kj}$.

Our goal is to determine the decompositions W and H . One way to do so is by minimizing the divergence penalty.

$$\mathcal{L}/D(X||WH) = \sum_{i,j} X_{ij} \log \frac{X_{ij}}{(WH)_{ij}} - X_{ij} + (WH)_{ij}$$

We solve the optimization problem by multiplicative update as opposed to additive updates. The update rules are:

$$W_{ik} \leftarrow W_{ik} \frac{\sum_j H_{kj} X_{ij} / (WH)_{ij}}{\sum_j H_{kj}}$$
$$H_{kj} \leftarrow H_{kj} \frac{\sum_i W_{ik} X_{ij} / (WH)_{ij}}{\sum_i W_{ik}}$$

We perform the updates until \mathcal{L} is small.

Important note: When performing the updates, $0 / 0$ divisions may occur in some cases. To avoid complications, I have can add a small number to the denominator.

I have ran the code on data from New York Times. Each line in the csv file corresponds to a single document. It gives information about the index of a word and the number of times the word occurs in that document. It is written in the format **index : count**.