

fasttext是facebook开源的一个词向量与文本分类工具，在学术上没有太多创新点，好处是模型简单，训练速度非常快。简单尝试可以发现，用起来还是非常顺手的，做出来的结果也不错，可以达到上线使用的标准。

简单说来，fastText做的事情，就是把文档中所有词通过lookup table变成向量，取平均之后直接用线性分类器得到分类结果。fastText和ACL-15上的deep averaging network(DAN)比较相似，是一个简化的版本，去掉了中间的隐层。论文指出了对一些简单的分类任务，没有必要使用太复杂的网络结构就可以取得差不多的结果。

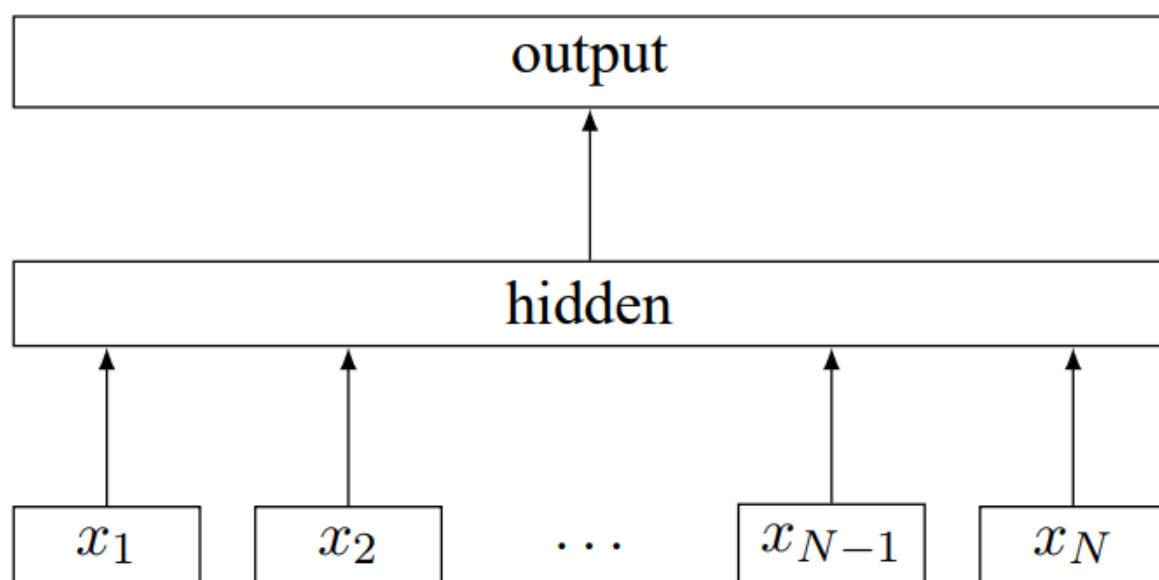


Figure 1: Model architecture of `fastText` for a sentence with N ngram features x_1, \dots, x_N . The features are embedded and averaged to form the hidden variable.

fastText论文中提到了一些tricks

- hierarchical softmax
 - 类别数较多时，通过构建一个霍夫曼编码树来加速softmax layer的计算，和之前word2vec中的trick相同
- N-gram features
 - 只用unigram的话会丢掉word order信息，所以通过加入N-gram features进行补充用hashing来减少N-gram的存储
- Subword
 - 对一些出现次数很少或者没有出现的词，使用subword的词向量之和来表达，如coresponse这个词，使用co的词向量与response的词向量之和来表示