

在工业界，很少直接将连续值作为逻辑回归模型的特征输入，而是将连续特征离散化为一系列0、1特征交给逻辑回归模型，优点有以下几个：

1. 离散特征的增加和减少都很容易，易于模型的快速迭代
2. 稀疏向量内积乘法运算速度快，计算结果方便存储，容易扩展
3. 离散化后的特征对异常数据有很强的鲁棒性，比如一个特征是年龄 $>30$ 是1，否则0。如果特征没有离散化，一个异常数据“年龄300岁”会给模型造成很大的干扰；
4. 逻辑回归属于广义线性模型，表达能力受限；单变量离散化为N个后，每个变量有单独的权重，相当于为模型引入了非线性，能够提升模型表达能力，加大拟合
5. 离散化后可以进行特征交叉，由M+N个变量变为M\*N个变量，进一步引入非线性，提升表达能力
6. 特征离散化后，模型会更稳定。比如如果对用户年龄离散化，20-30作为一个区间，不会因为一个用户年龄长了一岁就变成一个完全不同的人。当然处于区间相邻处的样本会刚好相反，所以怎么划分区间是门学问。
7. 特征离散化后，起到了简化逻辑回归模型的作用，降低了模型过拟合的风险。