

基本特征		统计特征		user answer		question answer			
用户id(1931654)	1	用户id_count	float64	1	感谢数 (最大、平均和求和)	1	问题关联回答数	1	0
问题id(1829900)	1	问题id_count	float64	1	点赞数 (最大、平均和求和)	1	问题上次回答时间 (day)	1	0
性别	1	性别_count	float64	1	收藏数 (最大、平均和求和)	1	问题最近一周回答比例	1	0
估值 (95-890)	1	访问频率	count	1	取赞数 (最大、平均和求和)	1			
访问频率	1	邀请创建时间-day(3838-3867,3868-3874)	1	1	举报数 (最大、平均和求和)	1			
邀请创建时间-day(3838-3867,3868-3874)	1	邀请创建时间-hour(0-23)	1	1	没有帮助数 (最大、平均和求和)	1			
问题创建时间-day(3838-3867,3868-3874)	1	问题创建时间-hour(0-23)	1	1	评论数 (最大、平均和求和)	1			
用户二分分类特征a	1	用户二分分类特征a_count	float64	1	被推荐次数 (最大、平均和求和)	1			
用户二分分类特征b	1	用户二分分类特征b_count	float64	1	用户的回答被标伏次数 (最大、平均和求和)	1			
用户二分分类特征c	1	用户二分分类特征c_count	float64	1	用户的回答被收入圈差次数 (最大、平均和求和)	1			
用户二分分类特征d	1	用户二分分类特征d_count	float64	1	是否有视频 (最大、平均和求和)	1			
用户二分分类特征e	1	用户二分分类特征e_count	float64	1	是否有图片 (最大、平均和求和)	1			
用户多分类特征a	1	用户多分类特征a_count	float64	1	内容字数 (最大、平均和求和)	1			
用户多分类特征b	1	用户多分类特征b_count	float64	1	用户回答时间与邀请时间距离-day (最小/最大/平均、均值、方差、总次数)	1			
用户多分类特征c	1	用户多分类特征c_count	float64	1	用户回答时间与邀请时间距离-hour (最小/最大/平均、均值、方差、总次数)	1			
用户多分类特征d	1	用户多分类特征d_count	float64	1	用户对回答的问题的感谢数统计 (出现最多的)	0			
用户多分类特征e	1	用户多分类特征e_count	float64	1	用户对回答过的话题embedding(average)	0			
邀请创建时间-weekday (1-7)	1	邀请创建时间-weekday_count		1	用户习惯回答时间-wkday	1			
问题创建时间-weekday (1-7)	1	问题创建时间-weekday_count		0	用户习惯回答时间-hour	1			
用户感兴趣/关注的话题文本embedding(average)	0	邀请创建时间-Friday Sat_Sun计数	0	0	用户关联回答数 (最近三天)	1			
问题相关的话题文本embedding(average)	0	问题创建时间与邀请时间距离-day (最小/最大/平均、均值、方差、总次数)	1	1	用户最近回答的数量 (最近三天)	1			
是否回答(10628265)	1	问题创建时间与邀请时间距离-hour (最小/最大/平均、均值、方差、总次数)	1	1	用户最近回答的数量 (最近一周)	1			
		用户拒绝邀请次数	0	0	用户最近回答的数量 (最近两周)	1			
		用户最感兴趣的topic	1	1	回答创建时间-day (3807-3867)	0			
		用户兴趣值的topic (最大、最小、均值、方差)	1	1					
		用户关注的topic的个数	1	1					
		用户感兴趣的话题_count	1	1					
		用户关注的话题_count	1	1					
		用户上次接受邀请距离当前时间-day	0	0					
		用户回答时间-hr_percent(用户id, hour)	0	0					
		用户回答时间-wk_percent(用户id, weekday)	0	0					
注意不能把label的信息置到特征中				颜色说明		用户回答话题数目的 (mean/std)			
				代表需要把时间线细分 (最近三天, 最近一周, 最近两周)		用户回答过话题的数目			
把用户根据回答历史情况拆开两个训练集, 面向两个训练模型				需要前缀分类对待		用户最擅长回答的话题id			
				代表证明有效的特征		用户回答最多的话题id			
				代表无效, 并且已经去除的特征					
按照目前向前链法切割训练集验证集 (交叉验证) 进行训练				代表注意事项		user question			
可以利用lightgbm分析特征的重要性						用户关注/和问题相关话题的重叠数目		1	
						用户感兴趣话题和问题相关话题的重叠数目		1	
						topic_intersection_values		0	
归一化一下下面出现的所有连续数值特征									
		lcc变换/区间放缩							
减少负样本, 用采样的方法									
		如何处理缺失值?							
		归一化的时候, 注意处理异常值							
		增加训练样本							
学习residual									
		用dnn去学习lab没学到的residual							
						</			