

用户预订售卖房型概率预测

团队介绍

队名：看命

成员：hwade

来自：华南理工大学

数据清洗

由于数据量比较大，数据无法一次加载如内存，所以在第一步读入数据时采用分块读入的方法，每次读入`chunksize`大小的数据，同时还要及时回收内存。

```
data = pd.read_csv(file_name, iterator=True)
loop = True
while loop:
    try:
        data_c = data.get_chunk(chunksize)
        yield data_c
        del data_c
        gc.collect()
    except StopIteration:
        loop = False
```

数据清洗

对一些没有区分度的特征则进行处理，能够进行转换的特征暂时留着，不能转换的特征直接删除。

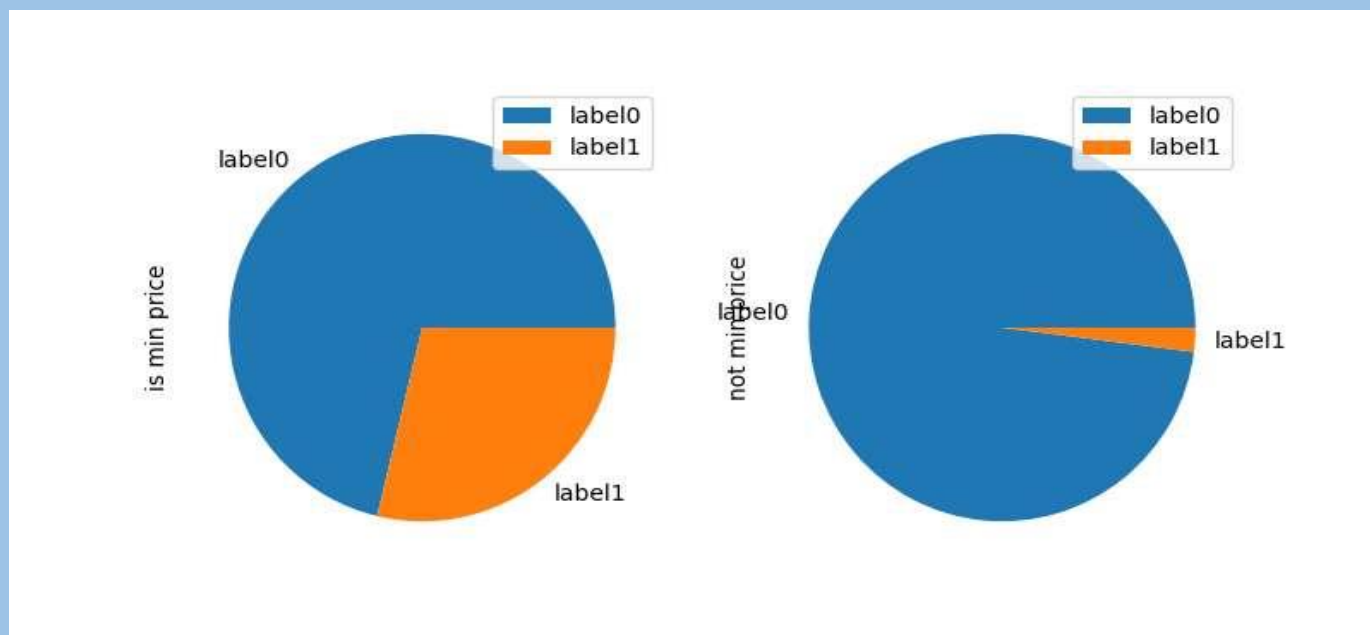
由于数据中存在大量的空值，如果简单的对含有空值的样本进行去除容易造成样本不足或者分布不平衡，所以采用不同方式对空值进行填充。

另外，需要采用低存储消耗的存储方式进行数据加载，如int64转成int32，float64转成float32。

```
for col in df_data.columns:
    if col.find('ratio') >= 0 or col.find('price') >= 0:
        df_data[col].fillna(df_data[col].mean(), inplace=True)
    elif col.find('roomservice') >= 0 or col.find('roomtag') >= 0:
        df_data[col].fillna(10, inplace=True)
    elif fill_method == 'mode':
        df_data[col].fillna(df_data[col].mode()[0], inplace=True)
    elif fill_method == 'mean':
        df_data[col].fillna(df_data[col].mean(), inplace=True)
    elif fill_method == 'quan':
        # quantile 0.5
        df_data[col].fillna(df_data[col].quantile(0.5), inplace=True)
    elif fill_method == 'pad':
        df_data[col].fillna(method='pad', inplace=True)
    elif fill_method == 'bfill':
        df_data[col].fillna(method='bfill', inplace=True)
```

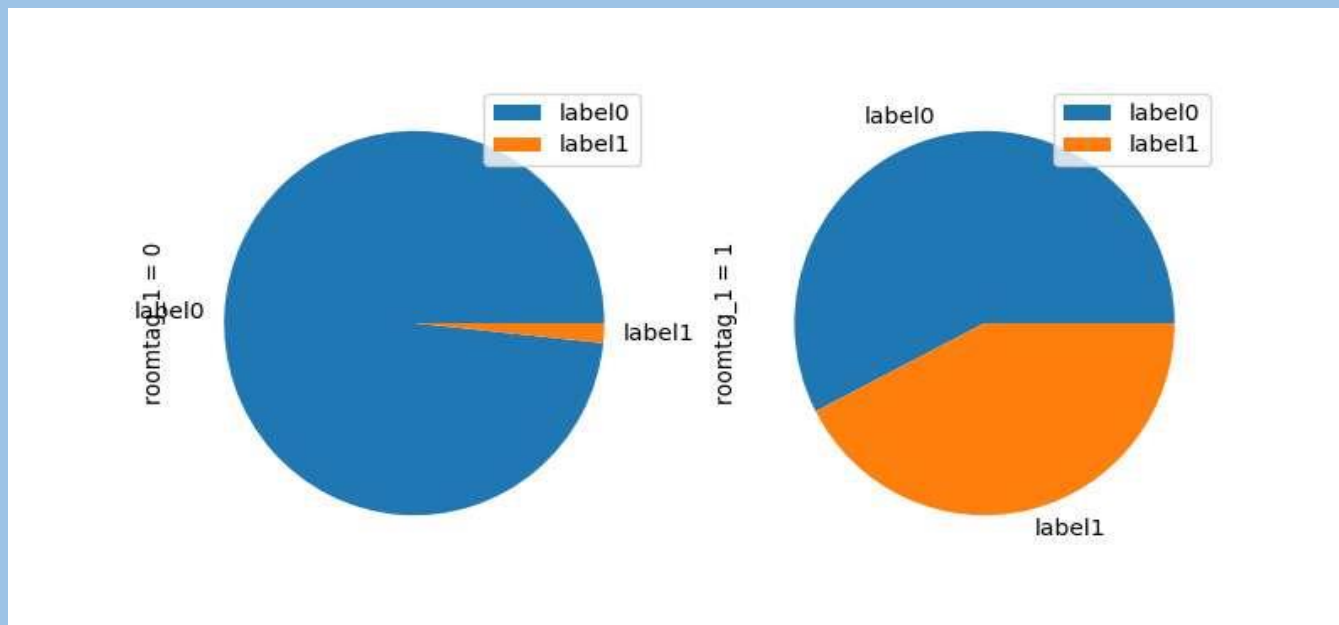
数据分析

数据中有很多维度的特征是针对价格的连续特征，从实际业务出发来思考，用户可能会因为价格的高低来决定是否下单，同一订单不同的房型中越低的价格越有可能下单。下图为标示是否为同一订单中价格最低的房型对最终是否下单的占比，可以看出最低价房型下单率明显比非最低价房型要高。



数据分析

下图表示物理房型roomtag_1为0和1的下单比率，可以看出roomtag_1为1时下单率接近50%，而为0时仅有不到5%的下单率。说明该特征是个强特征，在模型预测时具有很强的区分性。还有其他特征分析过程不一一列举。



特征工程

根据上面的特征构造方法，可以构造出以下的特征：

price_diff_ordermin	price_deduct 和同一订单中最小房价的一阶差分
price_diff_ordermax	price_deduct 和同一订单中最大房价的一阶差分
price_diff_orderavg	price_deduct 和同一订单中平均房价的一阶差分
diff_deal_price	price_deduct和用户平均下单价格的一阶差分
diff_workday_price	price_deduct和用户工作日平均定单价格的一阶差分
diff_holiday_price	price_deduct和用户假日平均定单价格的一阶差分
price_diff_useravg	price_deduct和用户平均定单价格的一阶差分
diff_return_promotion	returnvalue和用户平均返现价格的一阶差分
diff_star_price	price_deduct和同一星级star的所有房型平均价格的一阶差分
diff_rank_price	price_deduct和同一排位rank的所有房型平均价格的一阶差分
...	其他特征

特征工程

统计上一次订单记录与本条记录的差别，可以构造出以下的特征：

is_last_roomid	该房型roomid是否是上一次下单的房型
is_last_basicroomid	该物理房型basicroomid是否是上一次下单的物理房型
is_last_rank	该位置rank是否是上一次下单的rank
is_last_star	该房型星级star是否是上一次下单的星级star
is_last_roomtag_2	该房型标签roomtag_2是否是上一次下单的房型标签
is_last_roomservice_2	该服务roomservice_2是否是上一次下单的服务类型
...	其他特征

统计触发下单率类型的特征，可以构造出以下的特征：

roomtag_1_tri	该roomtag_1触发下单的次数
roomservice_1tri	该roomservice_1触发下单的次数
roomtag_1_tri_ratio	该roomtag_1触发下单的比率
roomservice_1_tri_ratio	该roomservice_1触发下单的比率
...	其他特征

统计方法直接涉及到orderlabel的值，所以需要做时间的滑窗来提取特征，否则容易造成数据泄露，如某一天的订单需要用订单以前的统计结果进行预测。由于时间关系，并未做滑窗提取。

对强特征进行特征组合，构造的特征如下：

roomtag_service_1_1	roomtag_1*10 + roomservice_1 的简单合并
roomtag_1_rank	roomtag_1*10 + rank 的简单合并
...	其他特征

将强特征进行组合，可以更加充分地挖掘出特征的使用价值，例如 **roomtag_1** 这个特征，单独使用时具有很强的区分性，若同时考虑与其他特征组合时，可以更大限度地挖掘出深层规律。但是由于时间关系，并未做更多特征组合的尝试。

