

CDA数据分析师毕业答辩

第五组

成员：陈 涛 彭飞舞
肖 威 刘 娜

2015年11月25日星期三

数据分析的力量

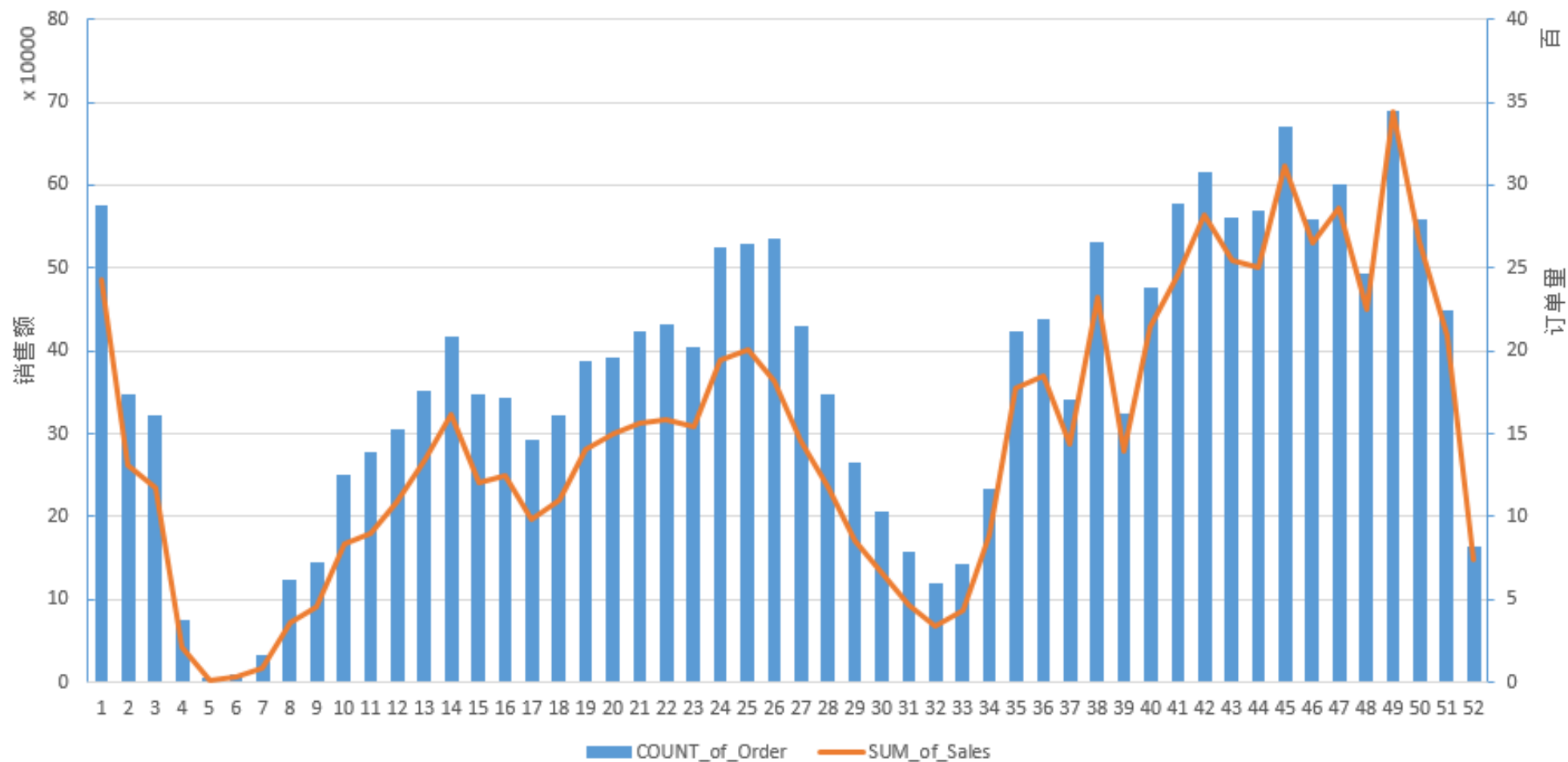
- ◆ 在购买不同的商品时，消费者如何分配他们的支出？
- ◆ 社会收入花费了多少，又节省了多少？
- ◆ 如何才能最好地测量与分析福利与贫穷？



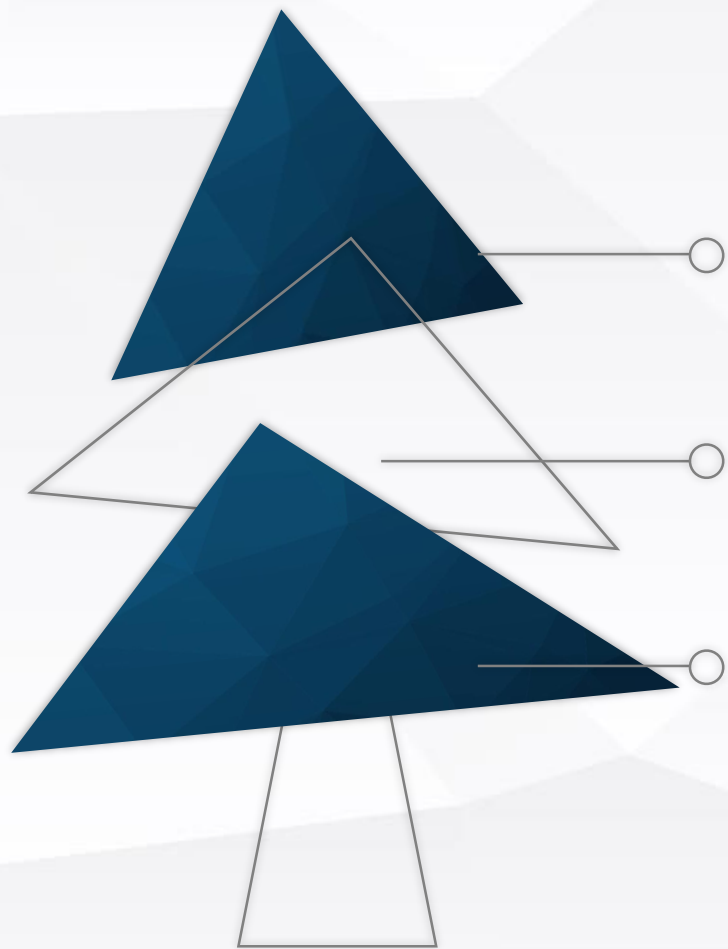
普林斯顿大学教授安格斯·迪顿 (Angus Deaton)
2015年诺贝尔经济学奖获得者

淘宝电商数据分析与挖掘

淘宝电商数据分析



目录



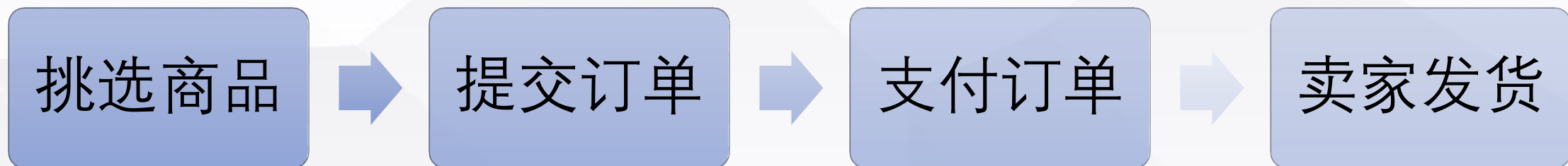
一、理解业务与数据

二、产品分析

三、客户分析

四、相关建议

理解业务与数据



CONTENTS PROCEDURE

订单数据

买家信息

商品信息

支付信息

其他信息

按字母排序的变量和属性列表

# 变量	类型	长度	输出格式	输入格式	标签
13 Address	字符	40	\$CHAR39.	\$CHAR39.	收货地址
23 City	字符	40			城市
3 Count_ID	字符	9	\$CHAR9.	\$CHAR9.	买家支付宝账号
24 Create_order_date	数值	8	DATE9.		订单创建日期
12 Message_num	数值	8	BEST12.	BEST12.	买家留言字数
2 Name	字符	11	\$CHAR11.	\$CHAR11.	买家会员名
1 Order	数值	8	BEST16.	COMMA12.	订单编号
20 Order_note	数值	8	BEST12.	BEST12.	订单备注
11 Order_status	字符	8	\$CHAR8.	\$CHAR8.	订单状态
9 Pay_amount	数值	8	BEST12.	BEST12.	买家实际支付金额
25 Pay_date	数值	8	DATE9.		订单支付日期
10 Pay_integration	数值	8	BEST12.	BEST12.	买家实际支付积分
16 Pay_time	数值	8	DATETIME18.	DATETIME18.	订单付款时间
17 Product_name	字符	2680	\$CHAR1782.	\$CHAR1782.	宝贝标题
21 Product_num	数值	8	BEST12.	BEST12.	宝贝总数量
18 Product_type	数值	8	BEST12.	BEST12.	宝贝种类
22 Province	字符	40			省份
4 Sales	数值	8	BEST12.	BEST12.	买家应付货款
7 Total_amount	数值	8	BEST12.	BEST12.	总金额

观测变量索引	120757
观测长度	36
删除的观测	0
已压缩	8416
已排序	0
	NO
	NO
(EUC)	

继续理解数据

甄别无用变量

- ◆ 支付积分
- ◆ 返点积分
- ◆ 实际支付积分
- ◆ 宝贝种类
- ◆

生成需要变量

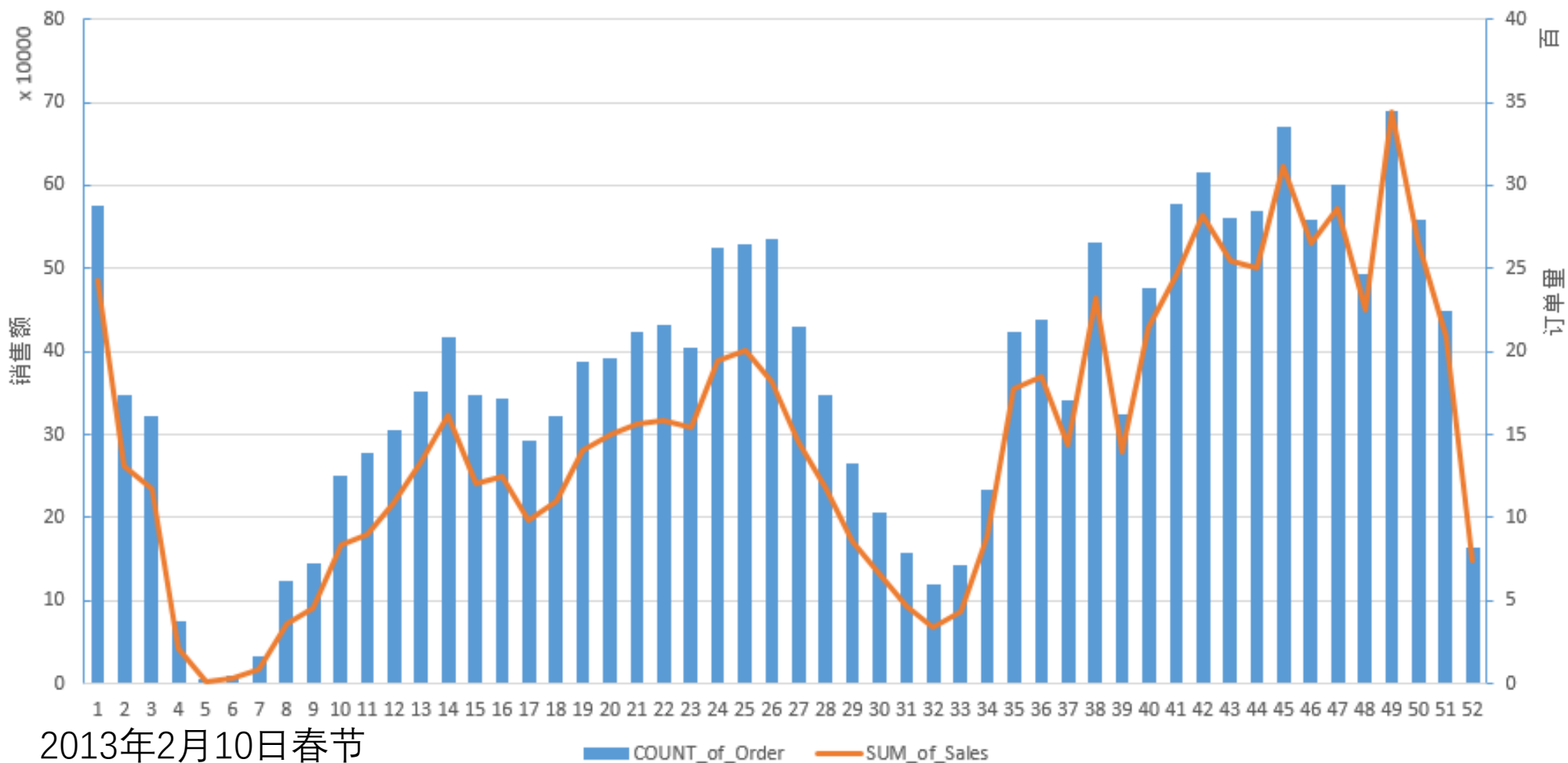
- ◆ 从收货地址中提取省份
- ◆ 买家支付账户个数
- ◆ 提交、支付订单的时间间隔
- ◆ 购买频率
- ◆

```
/*创建变量*/  
data bysj.bysj;  
  set bysj.bysj;  
  Province = scan(address,1," ");  
  City = scan(address,2," ");  
  Create_order_date = datepart(create_order_time);  
  Pay_date = datepart(Pay_time);  
  timed = pay_date - create_order_date;  
  timeh = (pay_time - create_order_time)/60;  
  create_order_hour = hour(create_order_time);  
  month = month(create_order_date);  
  week = int((Create_order_date - '05JAN2013'd)/7) + 1;  
  weekday = weekday(Create_order_date);  
  price = sales/product_num;
```

```
proc means  
  data=by sj.product min mean max p1 p5 p10 p25 p50 p75 p90 p95 p99;  
  var price;  
  class month;  
quit;
```

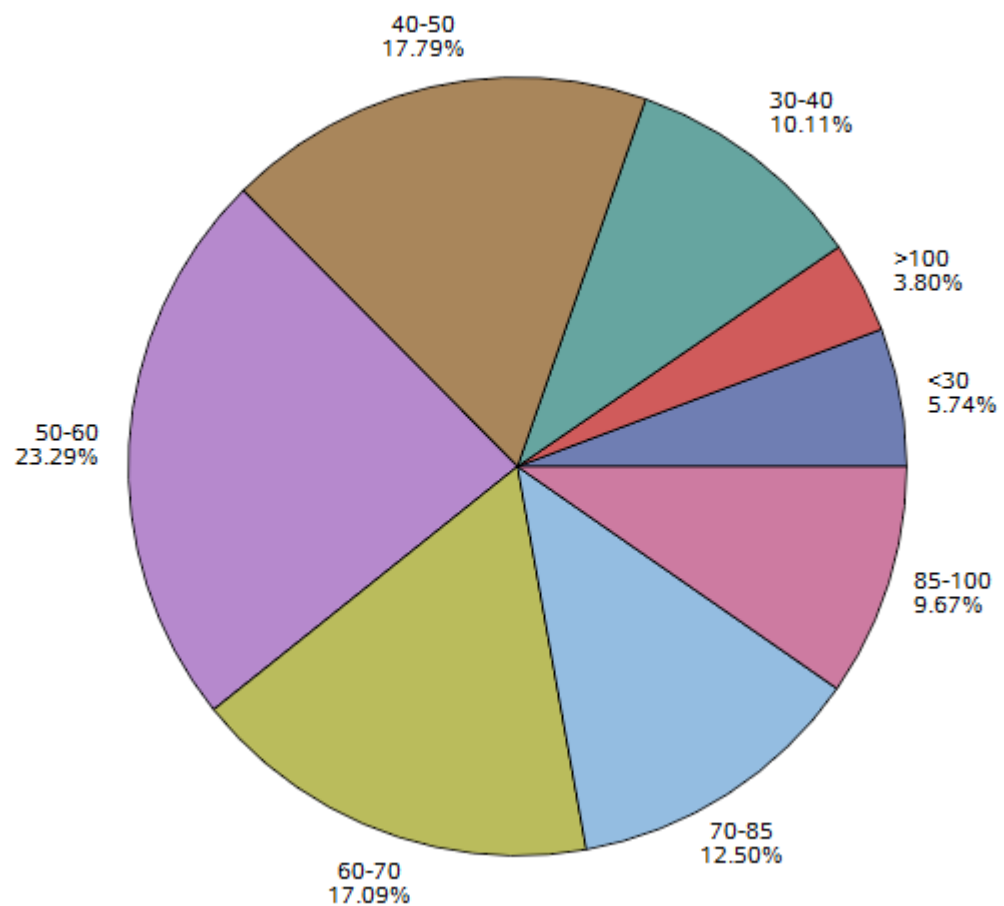
初步认识数据

淘宝电商数据分析



简单描述——商品分析

商品价格区间分布图



个百分位数	第 50 个百分位数	第 75 个百分位数	第 90 个百分位数	第 95 个百分位数	第 99 个百分位数
00000	66.3600000	82.7600000	100.5250000	139.0000000	228.0000000
00000	59.3400000	73.9750000	98.0000000	116.0500000	149.0000000
00000	56.0375000	70.4520000	89.0000000	95.0000000	114.0000000
00000	59.0000000	79.0000000	95.0000000	96.3333333	119.0000000

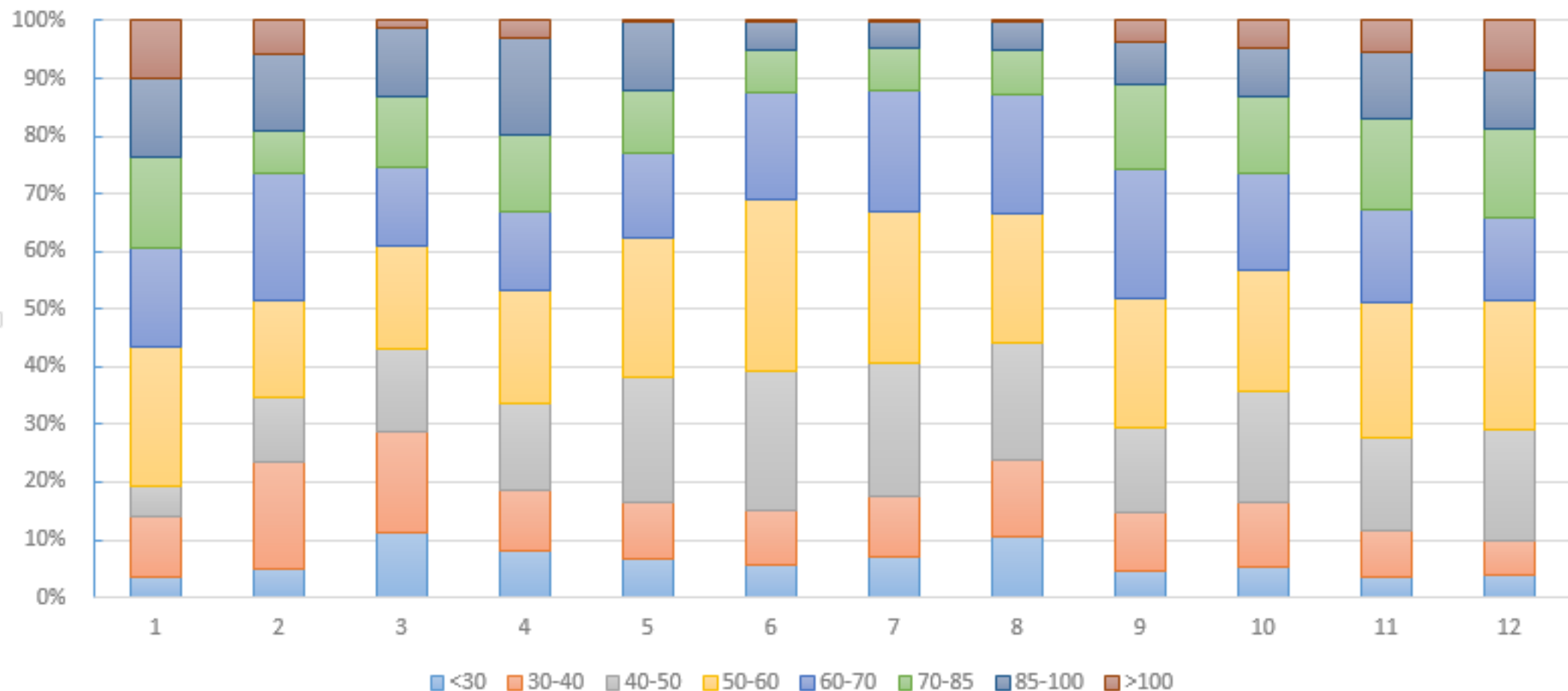
/*价格分段*/

```
data bysj.product;  
set bysj.product;  
length price_flag $ 12;  
if price < 30 then price_flag = "<30";  
else if price < 40 then price_flag = "30-40";  
else if price < 50 then price_flag = "40-50";  
else if price < 60 then price_flag = "50-60";  
else if price < 70 then price_flag = "60-70";  
else if price < 85 then price_flag = "70-85";  
else if price < 100 then price_flag = "85-100";  
else price_flag = ">100";  
format price_flag $8.;
```

run;

简单描述——商品分析

各价格区间商品销量占比

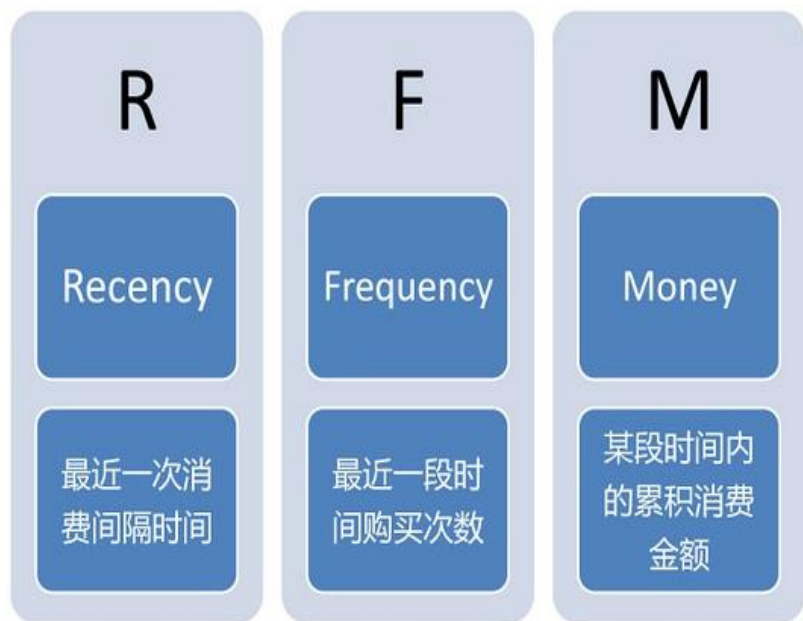


深入挖掘——客户分析

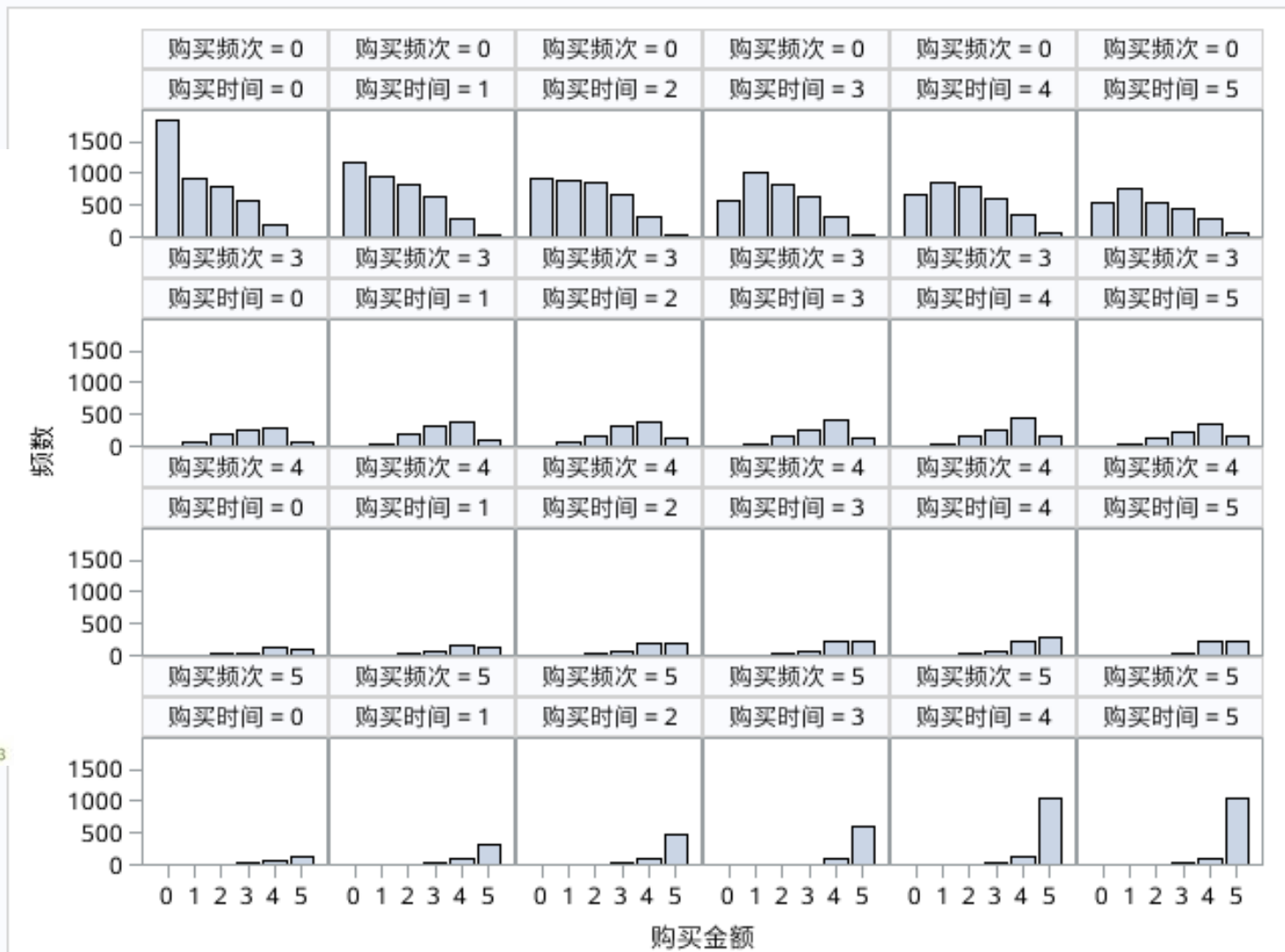
RFM客户分类模型

重复购买建模分析

RFM客户分类模型

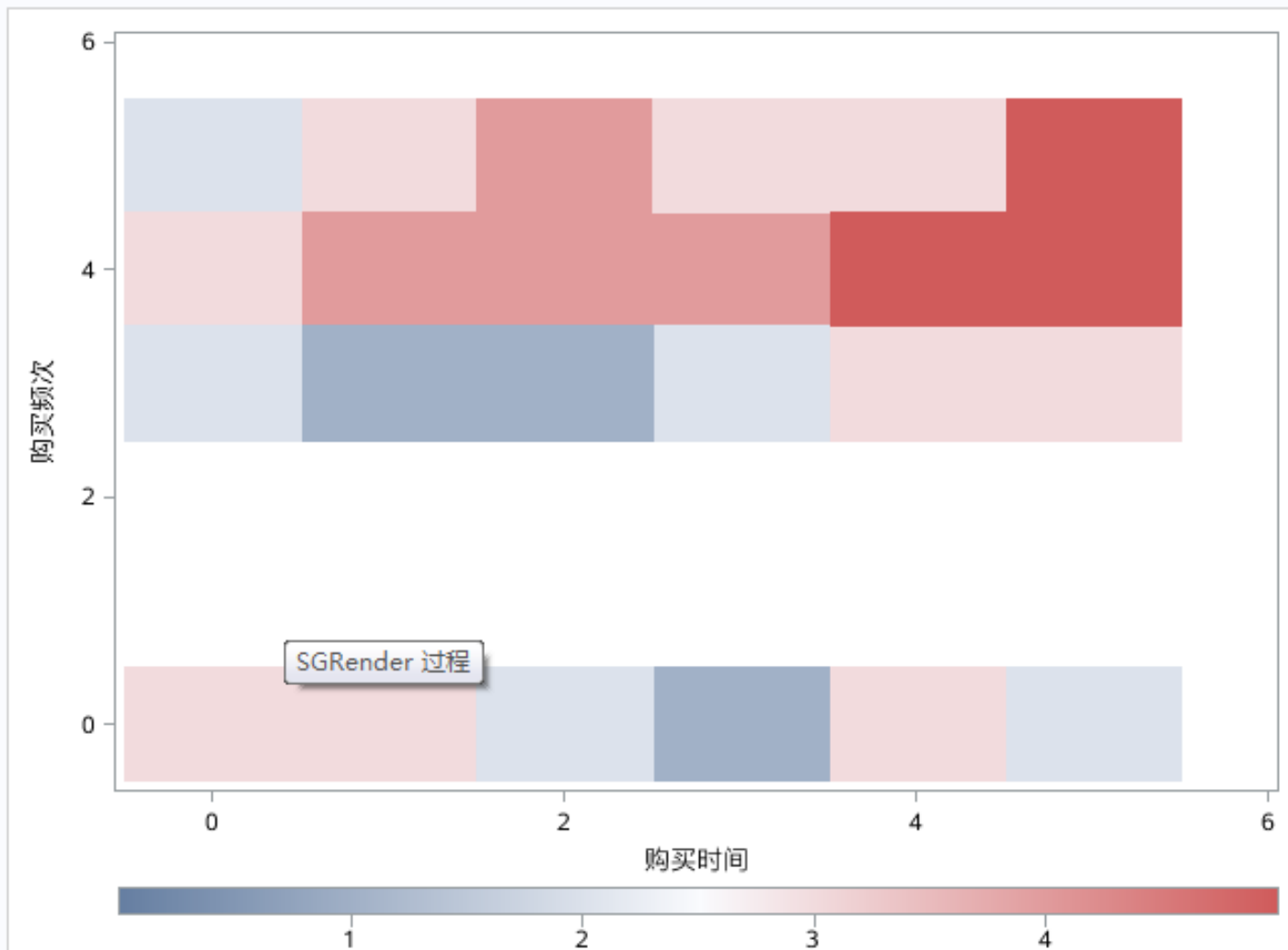


RFM得分分布



RFM客户分类模型

RFM得分热力图

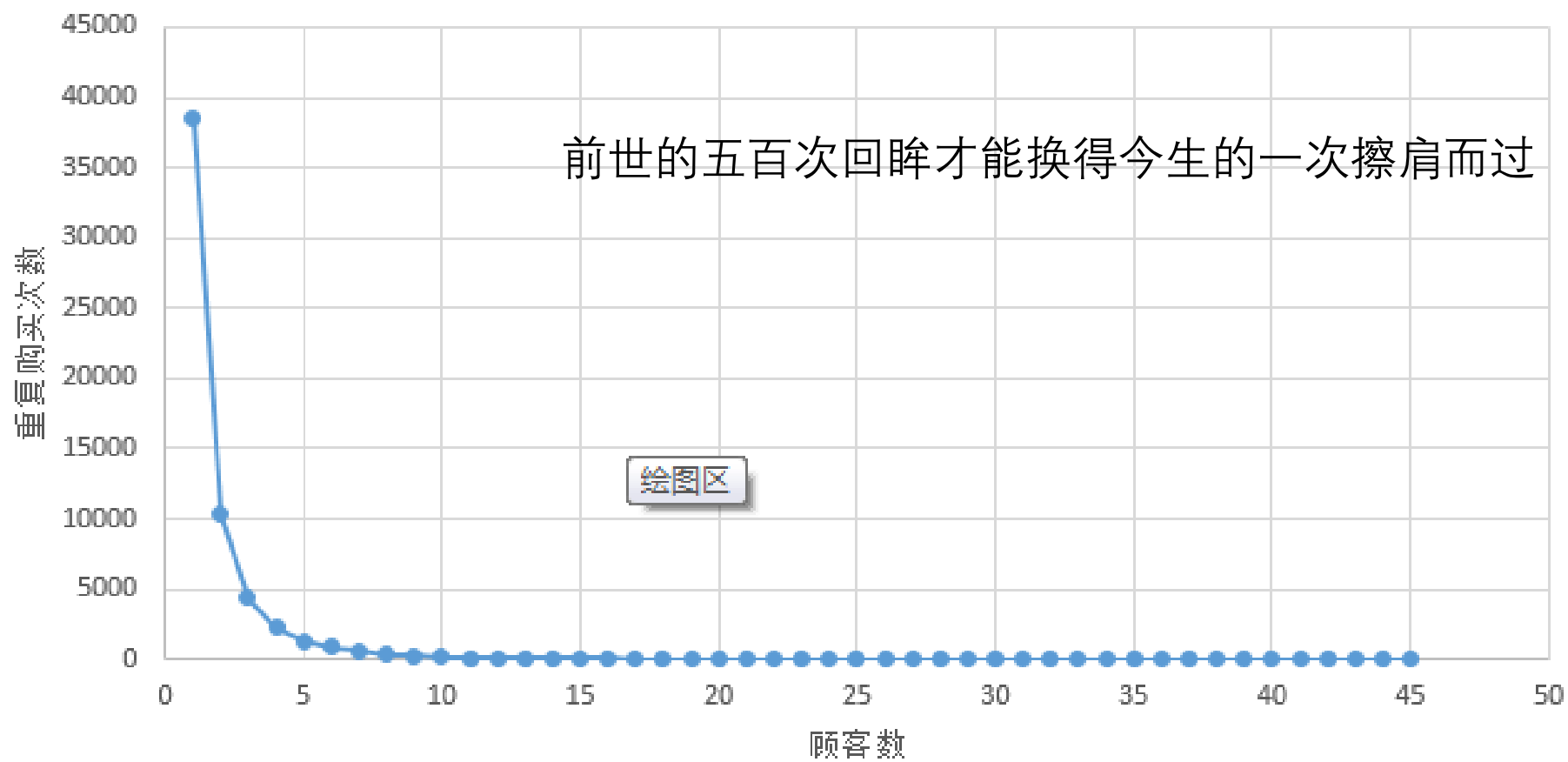


找准目标客户

对症下药

重复购买建模分析

重复购买次数分布图



重复购买建模分析

- ◆分割数据集：训练集70%，测试集30%
- ◆构建dmdb
- ◆Dmsplit方法筛选变量
- ◆建立logistic回归模型
- ◆测试模型
- ◆模型评价

The DMSPLIT Procedure

Effect Summary					
最大似然估计分析					
参数	自由度	估计	标准误差	Wald卡方	Pr > 卡方
预测概率和观测响应的关联					
Intercept					.01
trans_e		一致部分所占百分比	71.7	Somers D	0.441 .01
Count		不一致部分所占百分比	27.6	Gamma	0.444 .01
Total_a		结值百分比	0.7	Tau-a	0.211 .01
timeh		对	1058297778	c	0.721 .01
Pay_amount	1	0.000508	0.000208	5.9633	0.0146
timed	1	-0.2003	0.0693	8.3566	0.0038
timen			52		.69
Pay_amount			99		.91
timed			274		.73

重复购买建模分析

误分类表用以判断模型优劣

FREQ 过程

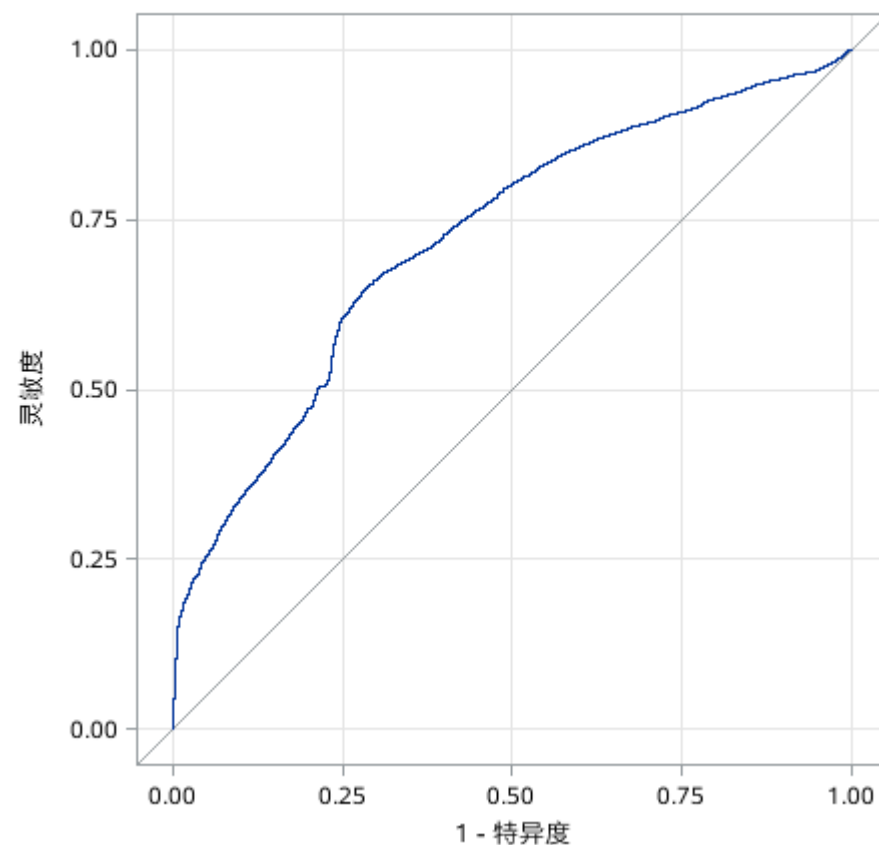
		I_rebuy(到: rebuy)		合计
		0	1	
F_rebuy(从: rebuy)	0	频数 6868	4301	11169
		百分比 24.19	15.15	39.34
	行百分比	61.49	38.51	
	列百分比	58.20	25.92	
1	频数	4933	12290	17223
		百分比 17.37	43.29	60.66
	行百分比	28.64	71.36	
	列百分比	41.80	74.08	
合计	频数	11801	16591	28392
	百分比	41.56	58.44	100.00

频数缺失 = 176

Pprob=0.45

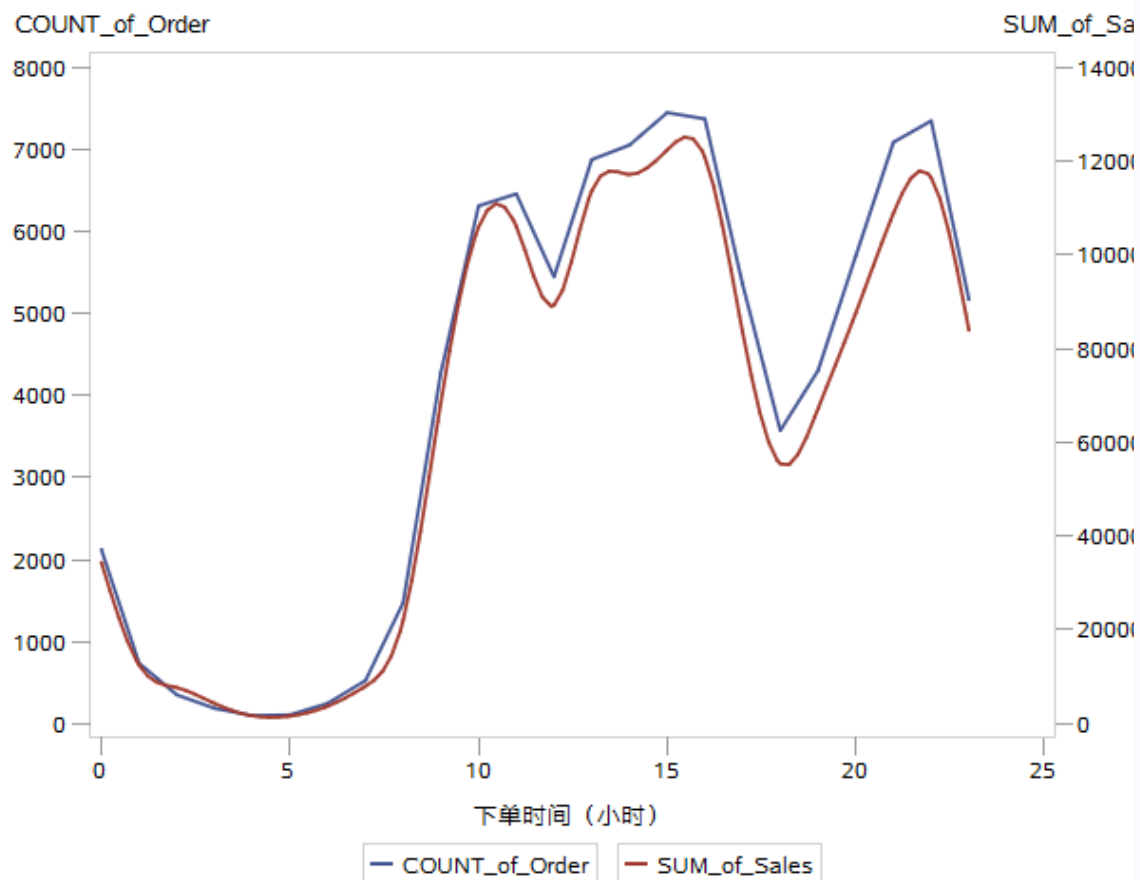
“模型”的 ROC 曲线

曲线下的面积= 0.7221

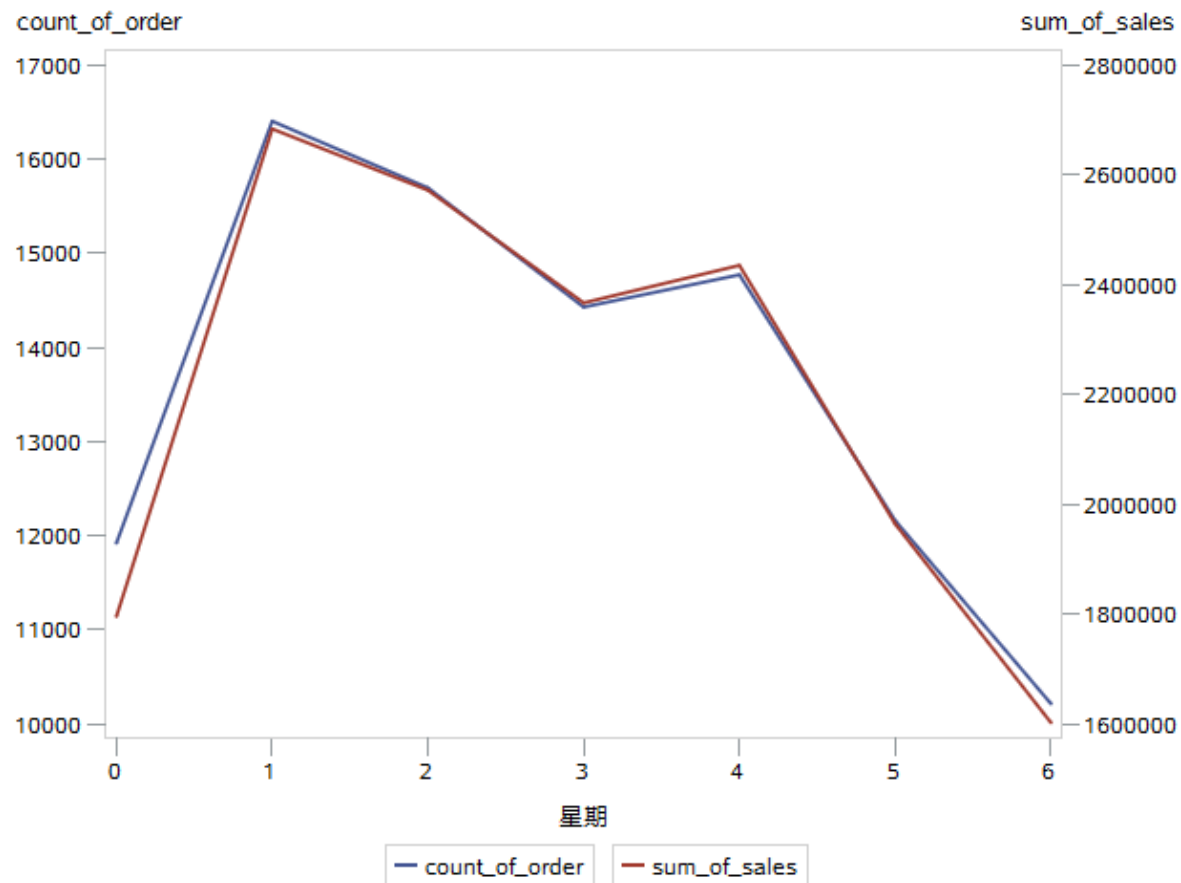


用户行为分析

顾客购买行为时间分布图

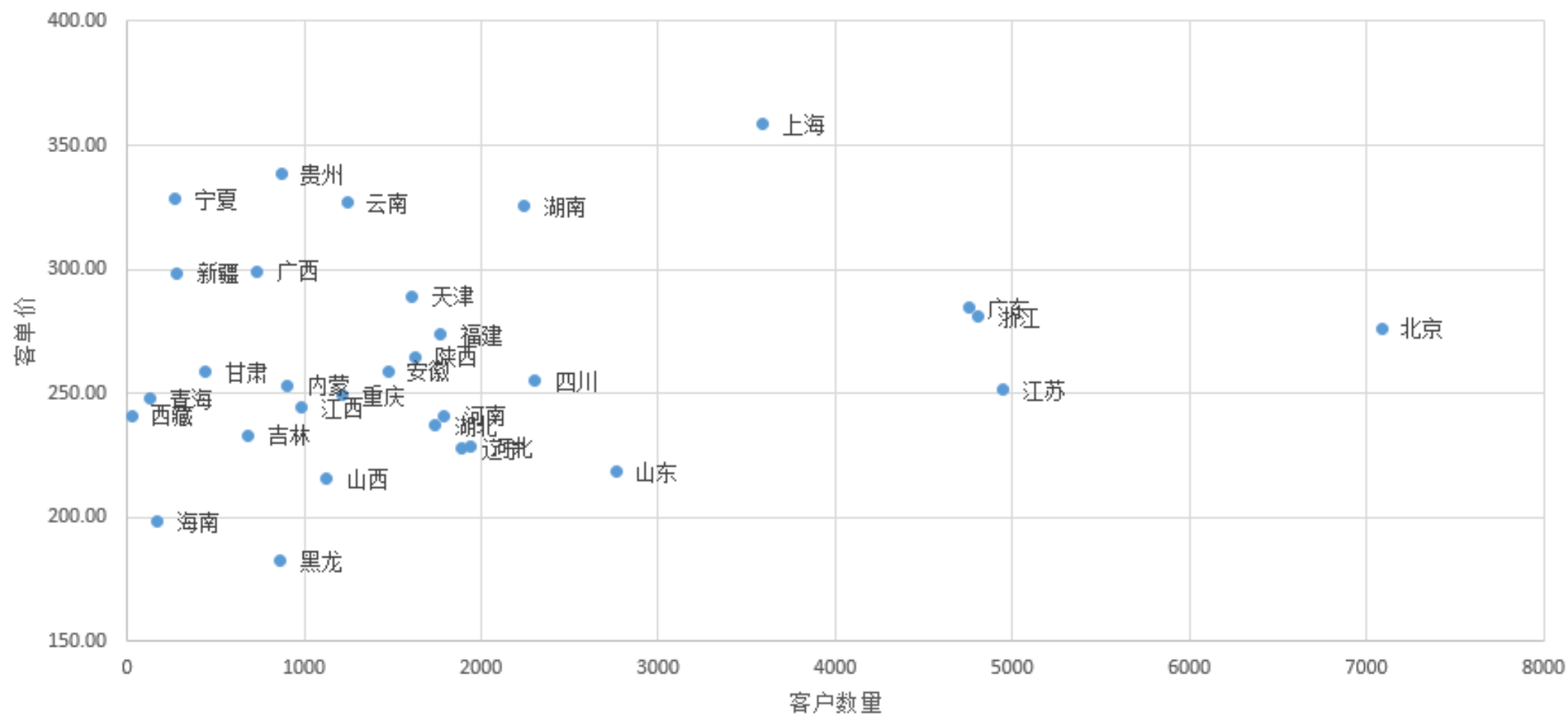


顾客购买行为分布图



用户行为分析

各省份客户数量与客单价关系图



相关建议

产品方面

客户关怀方面

客户服务方面

Q&A



THANKS