

# 预测商品册需求

## 执行概要

### 关键决策

公司今年是否向邮寄名单中新增的250名客户寄送产品目录册？

### 需要获取的数据

数据项	数据名称	数据来源	解释
1	Customer Segment	p1-customers.xlsx	在建模过程中建立虚拟变量
2	Avg Num Products Purchased	p1-customers.xlsx	在建模过程中建立预测变量
3	Avg Sale Amounts	p1-customers.xlsx	在建模过程中建立目标变量
4	Customer Segment	p1-mailinglist.xlsx	在分析过程中用作预测变量
5	Avg Num Products Purchased	p1-mailinglist.xlsx	在分析过程中用作预测变量
6	Avg Sale Amounts	p1-mailinglist.xlsx	在分析过程中用作目标变量
7	Score_Yes	p1-mailinglist.xlsx	顾客会对生产目录有所反应且进行购买的概率
8	平均毛利率50%	业务统计	通过产品目录册出售的所有产品的平均毛利率（价格减去成本）
9	成本6.5\$	业务统计	印刷和寄送每本产品目录册成本

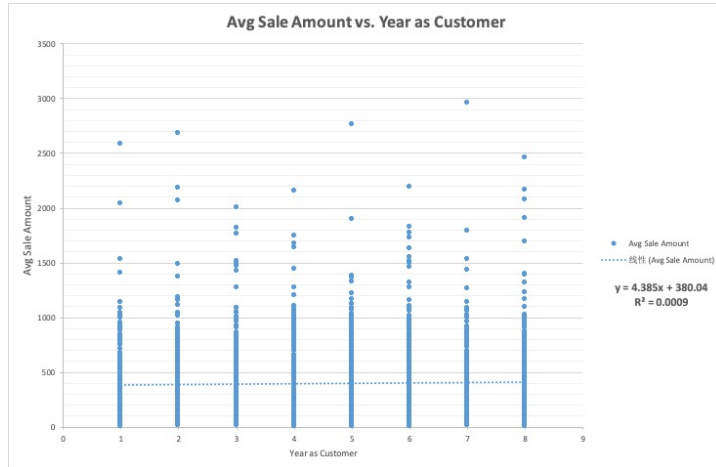
## 分析，建模和验证

1. 通过绘制每个数值型预测变量和目标变量之间的散点图来直观地查看是否有线性关系。

- Avg Num Products Purchased

随着平均购买产品数量的增加，平均销售额也以近似线性的趋势增长。





- Year as Customer

随着Year as Customer的增加，不会看到平均销售额增加（或减少）的趋势。可以得出结论，Year as Customer不是初步探索多元线性回归模型的一个很好的变量。

## 2. 对所有涉及到的单变量进行分析和回归:

数据项	数据分析对象	数据分析过程	数据分析结论																																																																																
1	Avg Num Products Purchased	<div><div>SUMMARY OUTPUT</div><div><div>回归统计</div><div>Multiple R 0.855754217 R Square 0.73231528 Adjusted R 0.732262476 标准误差 176.9070853 观测值 2375</div></div><div><div>方差分析</div><table><tr><th></th><th>df</th><th>SS</th><th>MS</th><th>F</th><th>Significance F</th></tr><tr><td>回归分析</td><td>1</td><td>201109435.1</td><td>201109435.1</td><td>6491.906448</td><td>0</td></tr><tr><td>残差</td><td>2373</td><td>73511948.03</td><td>30979.48632</td><td></td><td></td></tr><tr><td>总计</td><td>2374</td><td>274621383.1</td><td></td><td></td><td></td></tr></table><div><div>Coefficients</div><table><tr><th></th><th>标准误差</th><th>t Stat</th><th>P-value</th><th>Lower 95%</th><th>Upper 95%</th><th>下限 95.0%</th><th>上限 95.0%</th></tr><tr><td>Intercept</td><td>44.01516317</td><td>5.704322669</td><td>7.71610684</td><td>1.75315E-14</td><td>32.82919075</td><td>55.20113558</td><td>55.20113558</td></tr><tr><td>Avg Num Pro</td><td>106.2801833</td><td>1.319064914</td><td>80.57236777</td><td>0</td><td>103.6935443</td><td>108.8668224</td><td>108.8668224</td></tr></table></div></div></div> <div>R的平方=0.7323&gt;0.7, P&lt;0.05,具有统计显著性, 该数据可以参与多变量的线性回归分析</div>		df	SS	MS	F	Significance F	回归分析	1	201109435.1	201109435.1	6491.906448	0	残差	2373	73511948.03	30979.48632			总计	2374	274621383.1					标准误差	t Stat	P-value	Lower 95%	Upper 95%	下限 95.0%	上限 95.0%	Intercept	44.01516317	5.704322669	7.71610684	1.75315E-14	32.82919075	55.20113558	55.20113558	Avg Num Pro	106.2801833	1.319064914	80.57236777	0	103.6935443	108.8668224	108.8668224																																	
	df	SS	MS	F	Significance F																																																																														
回归分析	1	201109435.1	201109435.1	6491.906448	0																																																																														
残差	2373	73511948.03	30979.48632																																																																																
总计	2374	274621383.1																																																																																	
	标准误差	t Stat	P-value	Lower 95%	Upper 95%	下限 95.0%	上限 95.0%																																																																												
Intercept	44.01516317	5.704322669	7.71610684	1.75315E-14	32.82919075	55.20113558	55.20113558																																																																												
Avg Num Pro	106.2801833	1.319064914	80.57236777	0	103.6935443	108.8668224	108.8668224																																																																												
2	Customer Segment	<div><div>SUMMARY OUTPUT</div><div><div>回归统计</div><div>Multiple R 0.9148156 R Square 0.83688758 Adjusted R 0.83618768 标准误差 137.512507 观测值 2366</div></div><div><div>方差分析</div><table><tr><th></th><th>df</th><th>SS</th><th>MS</th><th>F</th><th>Significance F</th></tr><tr><td>回归分析</td><td>5</td><td>229065911</td><td>45813182.3</td><td>3028.419779</td><td>0</td></tr><tr><td>残差</td><td>2361</td><td>44645777</td><td>18909.6895</td><td></td><td></td></tr><tr><td>总计</td><td>2366</td><td>273711688</td><td></td><td></td><td></td></tr></table><div><div>Coefficients</div><table><tr><th></th><th>标准误差</th><th>t Stat</th><th>P-value</th><th>Lower 95%</th><th>Upper 95%</th><th>下限 95.0%</th><th>上限 95.0%</th></tr><tr><td>Intercept</td><td>585.799494</td><td>14.839192</td><td>39.4765089</td><td>3.6317E-262</td><td>556.700295</td><td>614.898693</td><td>614.898693</td></tr><tr><td>Credit Card</td><td>-281.51137</td><td>11.9423762</td><td>-23.572476</td><td>1.6461E-110</td><td>-304.93</td><td>-258.09274</td><td>-304.93</td></tr><tr><td>Loyalty Clu</td><td>0</td><td>0</td><td>65535</td><td>#N/M!</td><td>0</td><td>0</td><td>0</td></tr><tr><td>Loyalty Clu</td><td>-431.44536</td><td>12.7031437</td><td>-33.963668</td><td>#N/M!</td><td>-456.35583</td><td>-406.53488</td><td>-406.53488</td></tr><tr><td>Store Mail</td><td>-527.65274</td><td>13.8732322</td><td>-38.033872</td><td>2.5671E-247</td><td>-554.85773</td><td>-500.44776</td><td>-500.44776</td></tr><tr><td>Avg Num Pro</td><td>66.9080778</td><td>1.51779131</td><td>44.0825281</td><td>0</td><td>63.9317357</td><td>69.8844199</td><td>69.8844199</td></tr></table></div></div></div> <div>R的平方=0.8369&gt;0.7, P&lt;0.05,具有统计显著性, 该数据可以参与多变量的线性回归分析</div>		df	SS	MS	F	Significance F	回归分析	5	229065911	45813182.3	3028.419779	0	残差	2361	44645777	18909.6895			总计	2366	273711688					标准误差	t Stat	P-value	Lower 95%	Upper 95%	下限 95.0%	上限 95.0%	Intercept	585.799494	14.839192	39.4765089	3.6317E-262	556.700295	614.898693	614.898693	Credit Card	-281.51137	11.9423762	-23.572476	1.6461E-110	-304.93	-258.09274	-304.93	Loyalty Clu	0	0	65535	#N/M!	0	0	0	Loyalty Clu	-431.44536	12.7031437	-33.963668	#N/M!	-456.35583	-406.53488	-406.53488	Store Mail	-527.65274	13.8732322	-38.033872	2.5671E-247	-554.85773	-500.44776	-500.44776	Avg Num Pro	66.9080778	1.51779131	44.0825281	0	63.9317357	69.8844199	69.8844199	
	df	SS	MS	F	Significance F																																																																														
回归分析	5	229065911	45813182.3	3028.419779	0																																																																														
残差	2361	44645777	18909.6895																																																																																
总计	2366	273711688																																																																																	
	标准误差	t Stat	P-value	Lower 95%	Upper 95%	下限 95.0%	上限 95.0%																																																																												
Intercept	585.799494	14.839192	39.4765089	3.6317E-262	556.700295	614.898693	614.898693																																																																												
Credit Card	-281.51137	11.9423762	-23.572476	1.6461E-110	-304.93	-258.09274	-304.93																																																																												
Loyalty Clu	0	0	65535	#N/M!	0	0	0																																																																												
Loyalty Clu	-431.44536	12.7031437	-33.963668	#N/M!	-456.35583	-406.53488	-406.53488																																																																												
Store Mail	-527.65274	13.8732322	-38.033872	2.5671E-247	-554.85773	-500.44776	-500.44776																																																																												
Avg Num Pro	66.9080778	1.51779131	44.0825281	0	63.9317357	69.8844199	69.8844199																																																																												
3	Responded to Last Catalog	<div><div>SUMMARY OUTPUT</div><div><div>回归统计</div><div>Multiple R 0.198931667 R Square 0.039573808 Adjusted R 0.039167536 标准误差 333.46859 观测值 2366</div></div><div><div>方差分析</div><table><tr><th></th><th>df</th><th>SS</th><th>MS</th><th>F</th><th>Significance F</th></tr><tr><td>回归分析</td><td>1</td><td>10831813.83</td><td>10831813.83</td><td>97.40725849</td><td>1.52919E-22</td></tr><tr><td>残差</td><td>2364</td><td>262879874.5</td><td>111201.3005</td><td></td><td></td></tr><tr><td>总计</td><td>2365</td><td>273711688.3</td><td></td><td></td><td></td></tr></table><div><div>Coefficients</div><table><tr><th></th><th>标准误差</th><th>t Stat</th><th>P-value</th><th>Lower 95%</th><th>Upper 95%</th><th>下限 95.0%</th><th>上限 95.0%</th></tr><tr><td>Intercept</td><td>417.7003098</td><td>7.117667086</td><td>58.68500236</td><td>0</td><td>403.7427925</td><td>431.6578271</td><td>431.6578271</td></tr><tr><td>If Responde</td><td>-261.301947</td><td>26.47567162</td><td>-9.86951156</td><td>1.52919E-22</td><td>-313.219892</td><td>-209.384003</td><td>-209.384003</td></tr></table></div></div></div> <div>R的平方=0.03960&lt;0.3, 不存在线性关系, 该数据不可以参与多变量的线性回归分析</div>		df	SS	MS	F	Significance F	回归分析	1	10831813.83	10831813.83	97.40725849	1.52919E-22	残差	2364	262879874.5	111201.3005			总计	2365	273711688.3					标准误差	t Stat	P-value	Lower 95%	Upper 95%	下限 95.0%	上限 95.0%	Intercept	417.7003098	7.117667086	58.68500236	0	403.7427925	431.6578271	431.6578271	If Responde	-261.301947	26.47567162	-9.86951156	1.52919E-22	-313.219892	-209.384003	-209.384003																																	
	df	SS	MS	F	Significance F																																																																														
回归分析	1	10831813.83	10831813.83	97.40725849	1.52919E-22																																																																														
残差	2364	262879874.5	111201.3005																																																																																
总计	2365	273711688.3																																																																																	
	标准误差	t Stat	P-value	Lower 95%	Upper 95%	下限 95.0%	上限 95.0%																																																																												
Intercept	417.7003098	7.117667086	58.68500236	0	403.7427925	431.6578271	431.6578271																																																																												
If Responde	-261.301947	26.47567162	-9.86951156	1.52919E-22	-313.219892	-209.384003	-209.384003																																																																												
4	Year as Customer	<div><div>SUMMARY OUTPUT</div><div><div>回归统计</div><div>Multiple R 0.029781864 R Square 0.000868959 Adjusted R 0.000465926 标准误差 340.0365645 观测值 2375</div></div><div><div>方差分析</div><table><tr><th></th><th>df</th><th>SS</th><th>MS</th><th>F</th><th>Significance F</th></tr><tr><td>回归分析</td><td>1</td><td>243578.0156</td><td>243578.0156</td><td>2.106623132</td><td>0.146794828</td></tr><tr><td>残差</td><td>2373</td><td>274377805.1</td><td>115624.8652</td><td></td><td></td></tr><tr><td>总计</td><td>2374</td><td>274621383.1</td><td></td><td></td><td></td></tr></table><div><div>Coefficients</div><table><tr><th></th><th>标准误差</th><th>t Stat</th><th>P-value</th><th>Lower 95%</th><th>Upper 95%</th><th>下限 95.0%</th><th>上限 95.0%</th></tr><tr><td>Intercept</td><td>380.0388359</td><td>15.28292813</td><td>24.86688628</td><td>1.6908E-121</td><td>350.0695612</td><td>410.0081105</td><td>410.0081105</td></tr><tr><td>Years as C</td><td>4.384997179</td><td>3.021175081</td><td>1.451421073</td><td>0.146794828</td><td>-1.53941893</td><td>10.30941329</td><td>10.30941329</td></tr></table></div></div></div> <div>R的平方=0.0009, 几乎接近0, 不存在线性关系, 该数据不可以参与多变量的线性回归分析</div>		df	SS	MS	F	Significance F	回归分析	1	243578.0156	243578.0156	2.106623132	0.146794828	残差	2373	274377805.1	115624.8652			总计	2374	274621383.1					标准误差	t Stat	P-value	Lower 95%	Upper 95%	下限 95.0%	上限 95.0%	Intercept	380.0388359	15.28292813	24.86688628	1.6908E-121	350.0695612	410.0081105	410.0081105	Years as C	4.384997179	3.021175081	1.451421073	0.146794828	-1.53941893	10.30941329	10.30941329																																	
	df	SS	MS	F	Significance F																																																																														
回归分析	1	243578.0156	243578.0156	2.106623132	0.146794828																																																																														
残差	2373	274377805.1	115624.8652																																																																																
总计	2374	274621383.1																																																																																	
	标准误差	t Stat	P-value	Lower 95%	Upper 95%	下限 95.0%	上限 95.0%																																																																												
Intercept	380.0388359	15.28292813	24.86688628	1.6908E-121	350.0695612	410.0081105	410.0081105																																																																												
Years as C	4.384997179	3.021175081	1.451421073	0.146794828	-1.53941893	10.30941329	10.30941329																																																																												

3. 调整的 R 平方为0.8366> (.70) ,是强相关模型。对于所选择的预测变量，多元线性回归模型产生的 p 值都低于 0.05，所以具有统计学意义。以下是多元线性回归结果：

SUMMARY OUTPUT									
回归统计									
Multiple R	0.9148102								
R Square	0.8368777								
Adjusted R Square	0.8366024								
标准误差	137.48321								
观测值	2375								
方差分析									
	df	SS	MS	F	Significance F				
回归分析	4	2.3E+08	5.7E+07	3039.744236	0				
残差	2370	44796869	18901.6						
总计	2374	2.75E+08							
Coefficients									
	标准误差	t Stat	P-value	Lower 95%	Upper 95%	下限 95.0%	上限 95.0%		
Intercept	303.46	10.57571	28.6944	1.1227E-155	282.7249	324.20208	282.7249	324.20208	
Loyalty Club and	281.84	11.90986	23.6643	2.5804E-111	258.4839	305.19358	258.4839	305.19358	
Loyalty Club Only	-149.36	8.972755	-16.645	6.34584E-59	-166.951	-131.7605	-166.951	-131.7605	
Store Mailing Lis	-245.42	9.767776	-25.125	1.0503E-123	-264.572	-226.2635	-264.572	-226.2635	
Avg Num Products	66.98	1.51504	44.2075	0	64.00526	69.947147	64.00526	69.947147	

- 根据提供的数据，最佳线性回归方程：  
$$Y = 303.46 + 66.98 * (\text{Avg. Num Products Purchased}) + 281.84 (\text{If Type:Loyalty Club and Credit Card}) - 149.36 (\text{If Type: Loyalty Club Only}) - 245.42 (\text{If Type: Store Mailing List}) + 0 (\text{If Type: Credit Card Only})$$

## 建议

公司应该向这250个客户寄送产品目录册。

### 预测利润

- 使用线性回归模型来确定每个新客户的平均销售额(Avg Sale Amounts)
  - 平均销售额=303.46 + 66.98\* (Avg. Num Products Purchased )+ 281.84 (If Type:Loyalty Club and Credit Card) – 149.36 (If Type: Loyalty Club Only) – 245.42 (If Type: Store Mailing List) + 0 (If Type: Credit Card Only)
- 接下来，计算每位客户的预测销售额
  - 预测销售额=平均销售额 X 购买商品的概率（Score\_Yes）X 毛利率 – 宣传册成本
- 合计这250名新客户的预测销售额，得出发送新的产品目录的预测利润
  - 预测利润=21987.96 \$ >10000 \$

因此，公司应该向这250名客户寄送产品目录册。