



Excel数据分析师突击——从入门到精通到项目实战

第3周

【声明】 本视频和幻灯片为炼数成金网络课程的教学资料，所有资料只能在课程内使用，不得在课程以外范围散播，违者将可能被追究法律和经济责任。

课程详情访问炼数成金培训网站

<http://edu.dataguru.cn>

关注炼数成金企业微信



■ 提供全面的数据价值资讯，涵盖商业智能与数据分析、大数据、企业信息化、数字化技术等，各种高性价比课程信息，赶紧掏出您的手机关注吧！

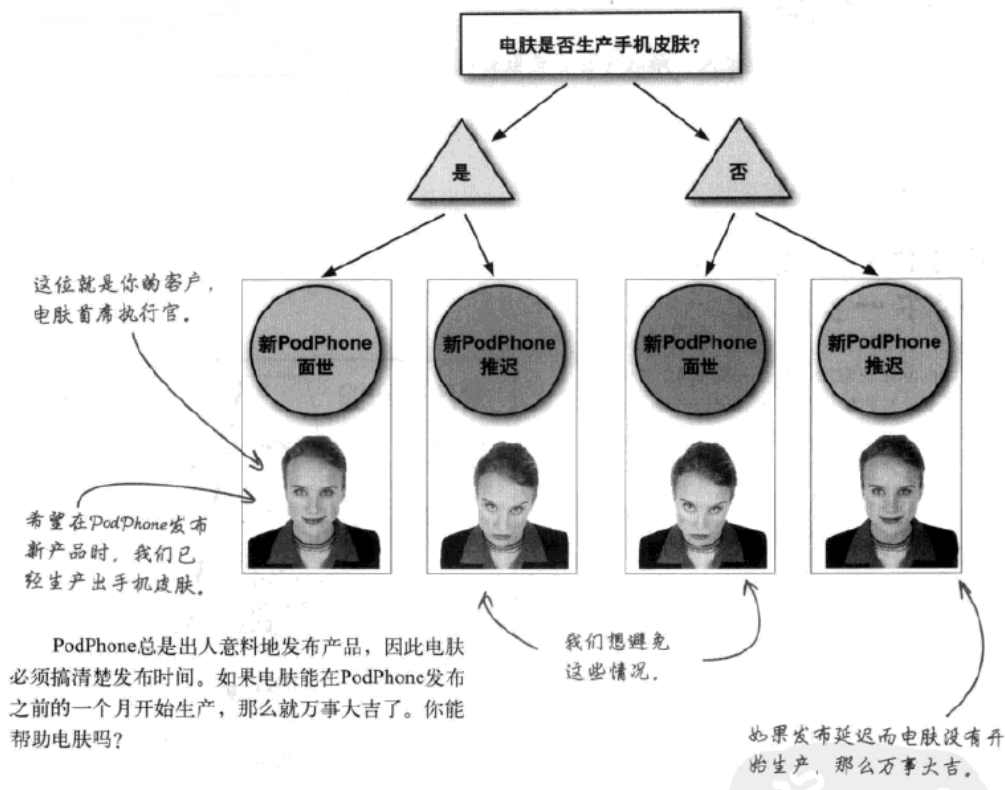


假设检验

世事纷纭，真假难辨。如何根据现有数据，推断出事实的真相？

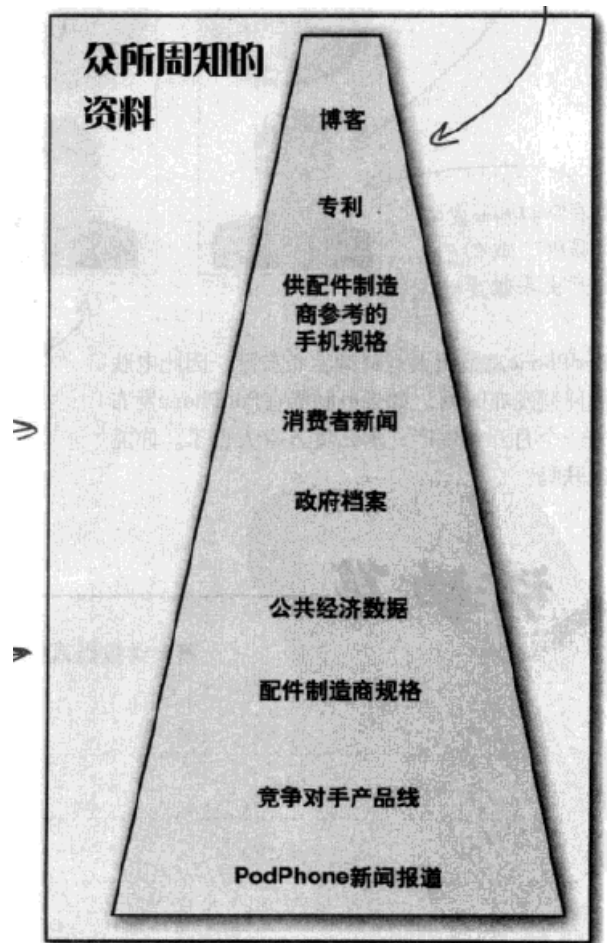
一个艰难的决定

- ◆ 电肤公司正在面临一个难题：是否需要将针对新PodPhone的手机皮肤生产提上日程？



一个艰难的决定

◆ 资料收集



PodPhone在新产品上的投资超过所有其他公司。

和竞争对手的手机相比，他们的手机性能将大幅改进。

PodPhone首席执行官说“我们绝不可能在明天推出新手机”。

一家竞争对手刚刚发布了一款性能优越的新手机。

经济回暖，消费者支出增多，正是卖手机的好时候。

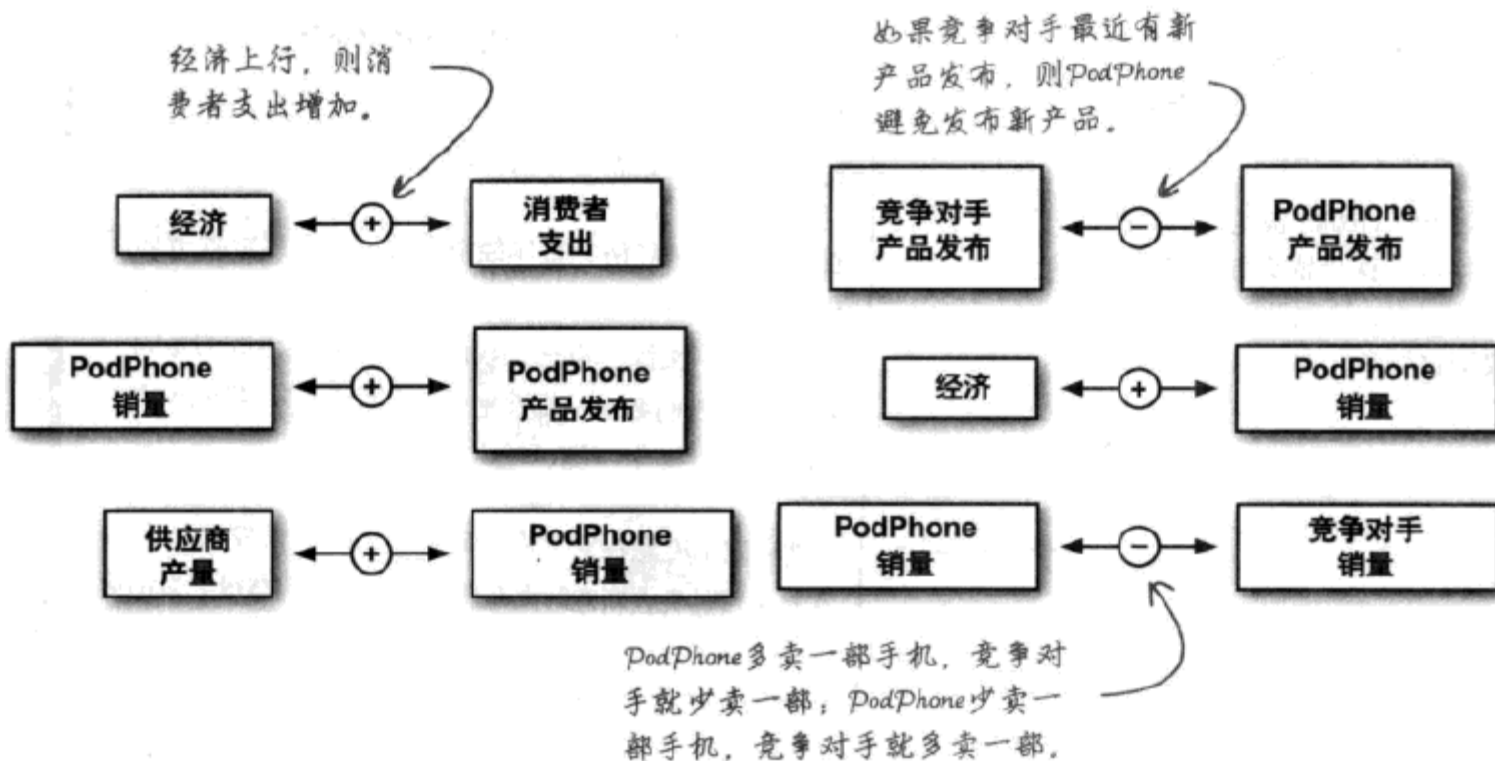
据传，PodPhone首席执行官表示一年以内不会发布新产品。

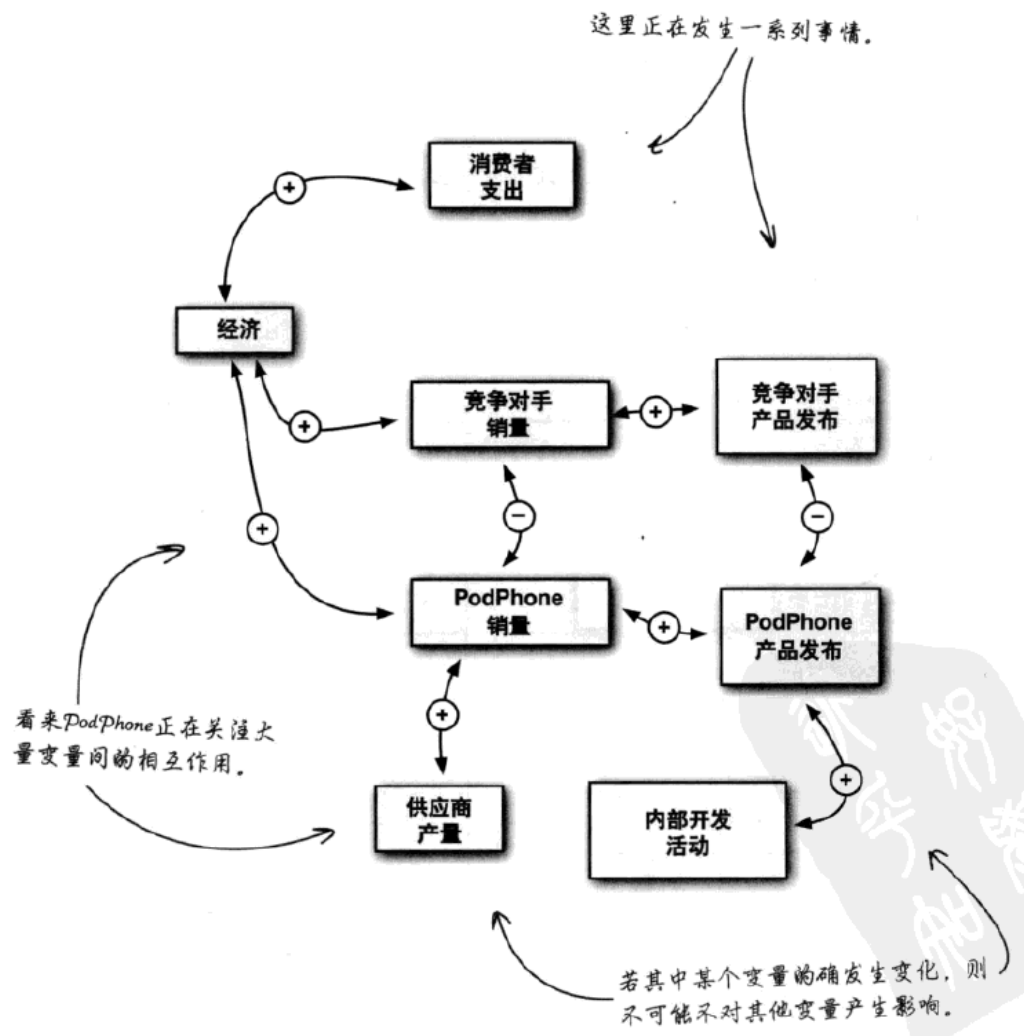
PodPhone手机发布战略备忘录

我们希望确定产品发布时间，以图实现最大销量，打败竞争对手，为此需要考虑种种因素。

首先关注的是经济，整体经济上行会促使消费者增加支出，经济下行则会抑制消费者支出，消费者支出是手机销量的唯一来源。但是，我们与竞争对手争夺的是同一块肥肉，我们多卖一部，他们就少卖一部；我们少卖一部，他们就多卖一部。

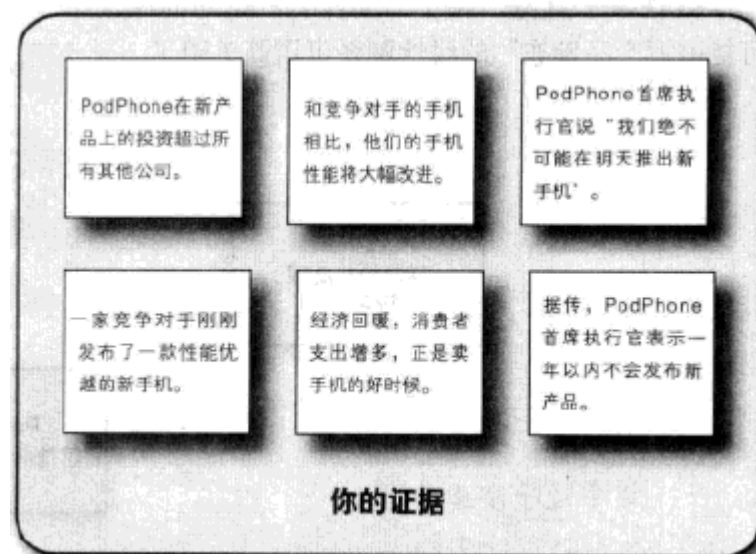
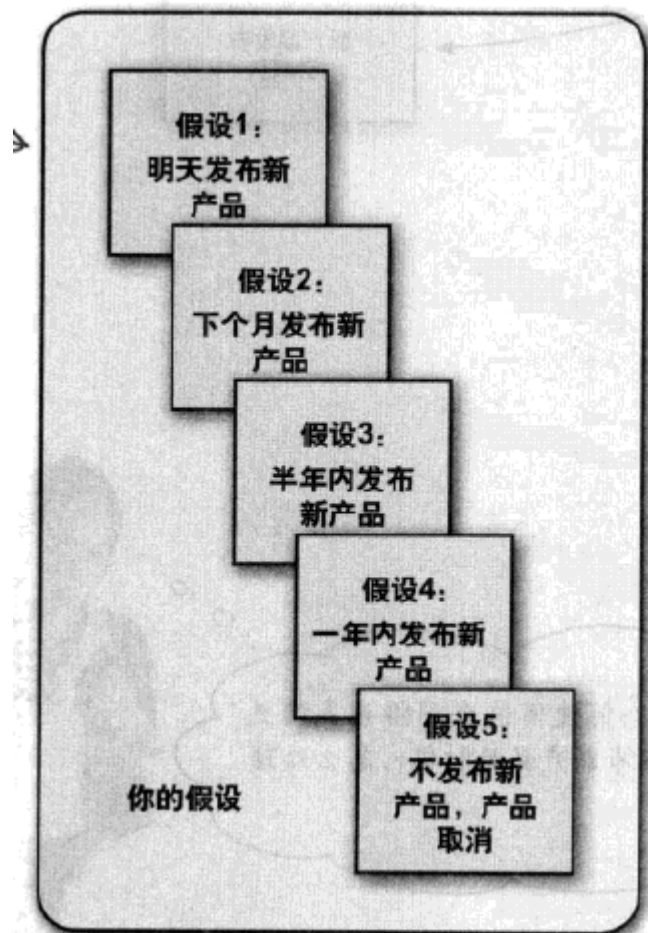
一般我们不愿意在对手有新手机上市的时候发布新产品，在对手产品失去新意时发布新产品会让我们夺得更多销量。我们的供应商和内部开发团队也限制了新手机生产能力。





◆ 提出几个假设

提交你的证据

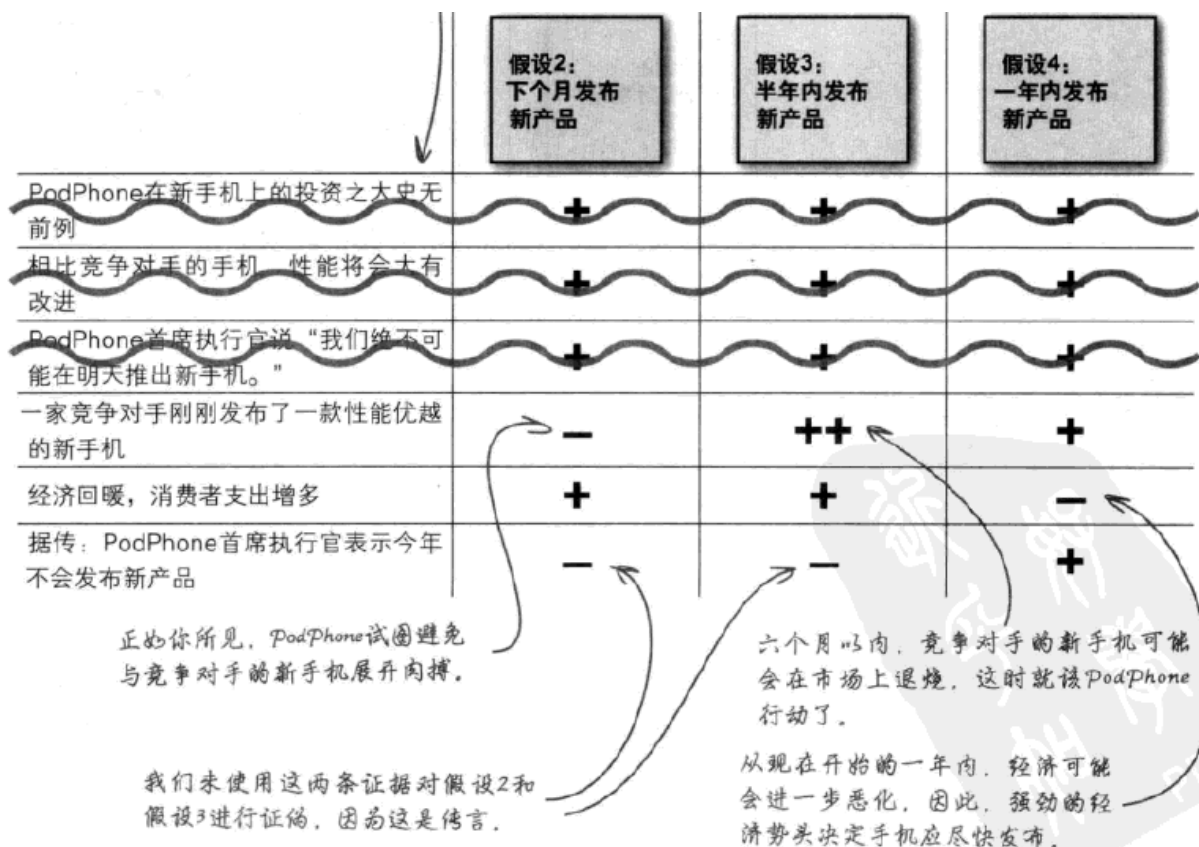


◆ 根据证据剔除错误假设



◆ 如何作出最佳推测？

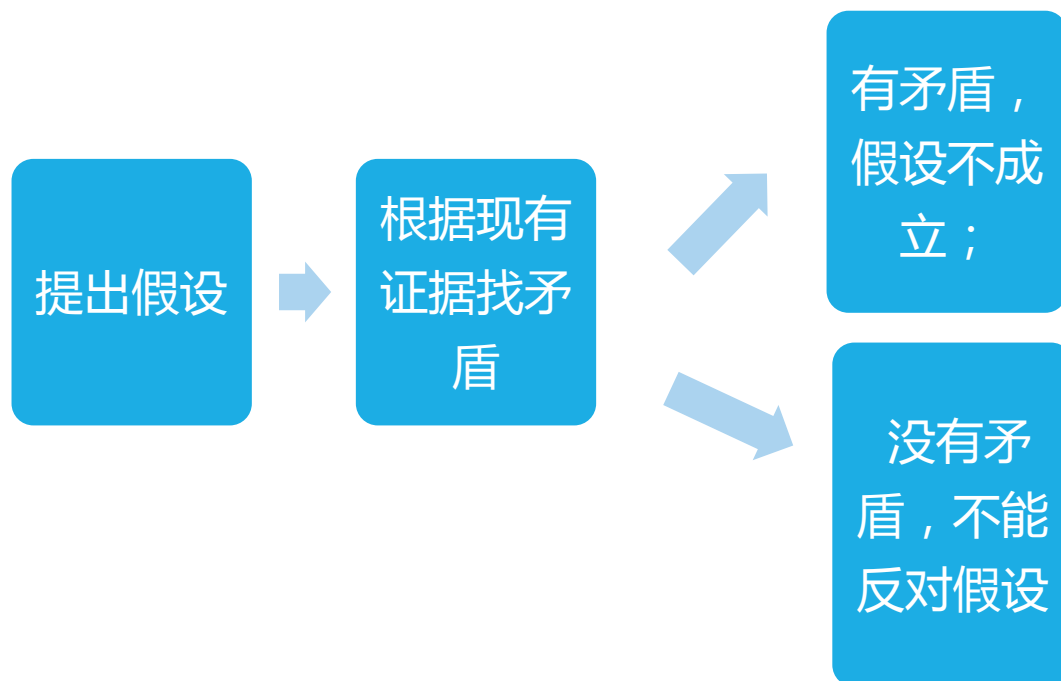
— 根据证据的诊断



◆ 添加新的证据

	假设2： 下个月发布 新产品	假设3： 半年内发布 新产品	假设4： 一年内发布 新产品
一家竞争对手刚刚发布了一款强大的新手机	—	++	+
经济回暖，消费者支出增多	+	+	—
据传：PodPhone首席执行官表示今年不会发布新产品	—	—	+
有人看见开发团队在开大型庆祝会，参加的人手里都拿着新手机。	+++	—	—

◆ 之前的例子中剔除错误假设的方法——伪证法的思路



◆ 提出假设——从检控官角度看

- 零假设/原假设——被告无罪
- 备择假设——被告有罪

◆ 什么是统计学上的零假设与备择假设？

◆ 零假设

◆ 备择假设

◆ 从证据中找矛盾

- 提出与被告无罪矛盾的证据
 - 犯罪现场的指纹，监控录像

◆ 统计学上的证据又是什么？——小概率事件

◆ 什么是小概率事件？

◆ 小概率事件在一次试验中出现说明了什么？

◆ 如何定义某个事件为小概率事件？发生的概率多小才算小概率事件？

- ◆ 得出结论
- ◆ 有小概率事件——与零假设有矛盾——拒绝零假设
- ◆ 没有小概率事件——没有发现与零假设的矛盾——不拒绝零假设
- ◆ 不拒绝 \neq 接受——可能只是证据不足
 - 疑点利益归于被告——被告被判无罪释放——被告是否真的无罪
 - 试验事件发生的概率大于设定阈值——不拒绝零假设

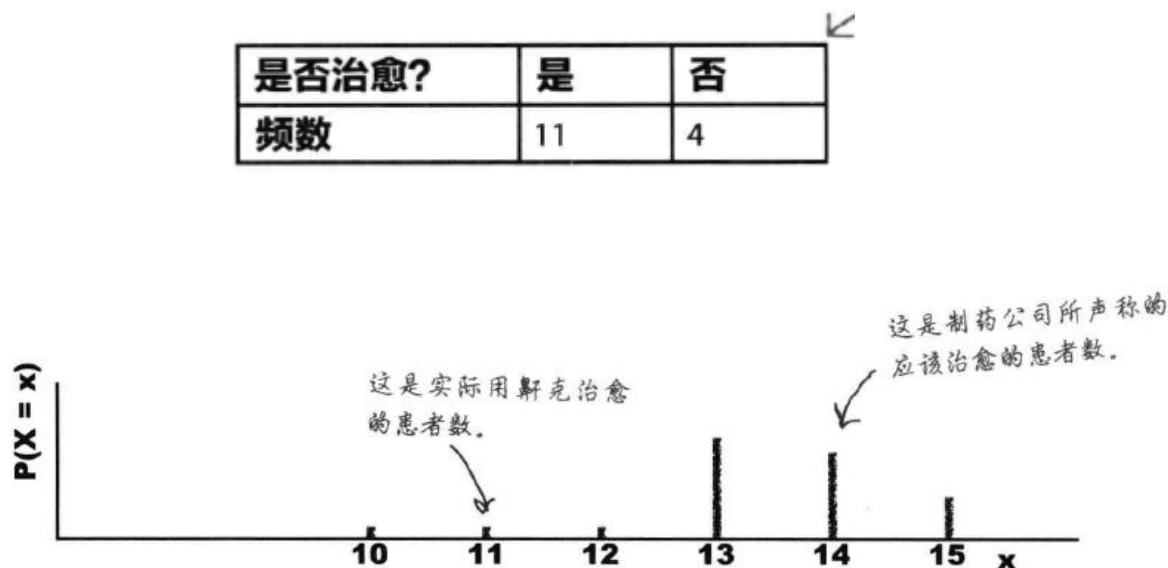
打鼾让你没精打采？
快让灵丹妙药“鼾克”来帮忙。
鼾克：患者2周内
治愈率90%。



新药鼾克，治打鼾有奇效！

治愈率很高的打鼾药

- ◆ 医生不相信这个治愈率，随机找了15位病人进行了2周的治疗，结果却不理想，问题出现在哪？



一起来做假设检验吧

即我们要对其进行试验的断言。

→ **1 确定要进行检验的假设**

2 选择检验统计量

← 我们需要选取能最有效地对断言进行检验的统计量。

我们需要使用某种确定性水平。

→ **3 确定用于做决策的拒绝域**

4 求出检验统计量的p值

← 我们需要了解在假定断言为真的情况下，我们的试验结果的可信程度。

5 查看样本结果是否位于拒绝域内

6 作出决策

← 接着需要了解试验结果是否位于确定性限值范围中。

一起来做假设检验吧

◆ 确定假设

- $H_0 : p=0.9$
- $H_1 : p<0.9$

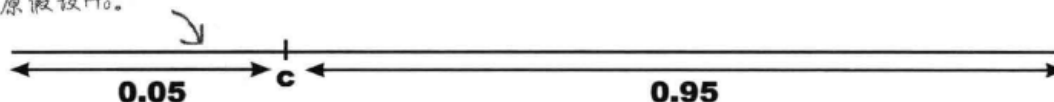
◆ 选择检验统计量

- 记治愈人数为 X , $X \sim B(15, 0.9)$

◆ 确定拒绝域——确定小概率事件的范围

◆ 计算P值——该次试验发生的概率

如果 $P(X \leq 11)$ 小于0.05, 说明数值11落在拒绝域中——我们可以拒绝原假设 H_0 。



◆ 作出决策

- 进行假设检验即选定一个断言，然后借助统计证据对其进行检验。
- 所检验的断言被称为原假设，用 H_0 表示。除非有有力的证据证明断言不正确，否则就接受断言。
- 备择假设即在有充分证据拒绝原假设 H_0 的情况下将接受的假设，用 H_1 表示。
- 检验统计量即用于对假设进行检验的统计量，是与检验具有最密切关系的统计量。选择检验统计量的时候，你假定 H_0 为真。
- 显著性水平用 α 表示，它表示你希望在观察结果的不可能程度达到多大时拒绝 H_0 。
- 拒绝域为一组数值，代表可用于否定原假设的最极端证据。选择拒绝域时，需考虑显著性水平，还要考虑用单尾还是双尾进行检验。
- 单尾检验的拒绝域位于数据的左侧或右侧，双尾检验的数据一分为二位于数距的两侧。可根据备择假设选择尾部。
- P值即取得样本结果或取得拒绝域方向上的更极端结果的概率。
- 如果P值位于拒绝域中，则有充足的理由拒绝原假设；如果P值位于拒绝域以外，则没有充足的证据。

- ◆ 样本量不足——证据不足，搜集新的证据再上诉

是否治愈?	是	否
频数	80	20

- ◆ 确定假设
- ◆ 计算检验统计量
- ◆ 确定拒绝域
- ◆ 计算P值
- ◆ 作出结论

假设检验需要证据。

进行假设检验时，你选取一个断言，然后对其进行试验。只有在有足够证据反驳这个断言时，你才能否定这个断言。这意味着检验是公正的，因为你做决策的唯一依据就是是否有充分证据。

如果我们一开始就接受医生的观点，就不会妥当地考虑证据。我们会在不考虑结果是否只能解释为偶然的情况下作出决策，而现在呢，我们有足够的证据表明，样本结果足以合理地拒绝原假设。这些结果具有统计显著性，因为它们不可能是偶然发生的。

◆ 错误类型

假设检验决策

实际情况	接受 H_0	拒绝 H_0
H_0 真	✓	第一类错误
H_0 假	第二类错误	✓

↑ ↑ ↑
这些都是错误 这给出检验的功效。

◆ 第一类错误 α

◆ 第二类错误 β

◆ 功效 $\text{Power}=1-\beta$

制药公司和他们的止咳糖浆制造厂发生了争议，厂方说注入药瓶的糖浆量符合正态分布 $X \sim N(355, 25)$ ，其中 X 是量得的每瓶糖浆容量，单位mL。制药公司用大样本进行了检验，发现100瓶糖浆的平均容量为356.5mL。请以1%的显著性水平检验厂方给出的均值假设，与此相对的另一说法是每瓶糖浆的容量均值大于355mL。

- ◆ **Dataguru（炼数成金）是专业数据分析网站，提供教育，媒体，内容，社区，出版，数据分析业务等服务。我们的课程采用新兴的互联网教育形式，独创地发展了逆向收费式网络培训课程模式。既继承传统教育重学习氛围，重竞争压力的特点，同时又发挥互联网的威力打破时空限制，把天南地北志同道合的朋友组织在一起交流学习，使到原先孤立的学习个体组合成有组织的探索力量。并且把原先动辄成千上万的学习成本，直线下降至百元范围，造福大众。我们的目标是：低成本传播高价值知识，构架中国第一的网上知识流转阵地。**
- ◆ **关于逆向收费式网络的详情，请看我们的培训网站 <http://edu.dataguru.cn>**

Thanks

FAQ时间