

1. 表 1.1 中若只包含编号为 1 和 4 的两个样例。
则给出相应的版本空间。

数据集:

色 洋 根 茎 高 声 如 例
1 青 绿 蜡 蜡 油 响 是
2 黄 黑 稍 蜡 油 响 否

假设空间

共 $3 \times 3 \times 3 \times 3 \times 3 \times 3$ 种假设。

下面我们逐一验证假设是否有

1. * * * * * 2. * 蜡 *
3. * * 油 * 4. 青 * *
5. * 蜡 油 6. 青 蜡 *
7. 青 * 油 8. 青 蜡 油。

得出出样例的结果。1.

∴ 版本空间只有 7 个，如上所示 2-8

2. 测试机器，学习能在互联网搜索的哪
些内容起作用。

国内搜索引擎：通过标签检索网页，集会对
网页进行分类标注。

推送服务：基于用户的点击，分析出用户的
潜在需求，从而推送广告。

自然语言处理：自动文本输入，自动识别自己。

3. 数据集 1000 个样例，正：反 = 1:1，划分为 20%
样例的训练集和 80% 的测试集，训练集出
有多少种划分方式。

对于测试集应有 150 个正例和 150 个反例。

∴ 共有 $\binom{150}{50}$ 种划分

4. 100 个样例，正：反 = 1:1，假设是学习算法所平均的模型误差。

将每个样例被误判为训练样例错误的类别，训练集用

10 折交叉验证法和留一法分别对模型进行评估的结果。

10 折交叉验证：因为每个子集都保持数据一般性，∴ 训练集中有
正例 = 反例 = 45。

测试集中随机抽取测试，留一法为 50%

留一法：训练集中正：反或反：正 = 49:50

对每种情况，被误判都错误，留一法为 100%

5. 训练集正例率 (TPR)，假设正例率 (FPR) 与查准率 (P)，
查全率 (R) 之间的关系。

	正	反
正	TP	FN
反	FP	TN

$P = \frac{TP}{TP + FP}$ 查准率 表示被正例为正的样例有多少
为真正的正样本。

$R = \frac{TP}{TP + FN}$ 查全率 表示所有正例中，有多少正例
被模型识别出。

两者一般是制约关系，不能同时增加。

$TPR = \frac{TP}{TP + FN}$ 正样本被识别数量与正例所占比例

$FPR = \frac{FP}{FP + TN}$ 被误判为正的负样本数量与负例所占比例

6. 线性正则化分析假设在线性所有数据上能获得最优解。
结果，对设计一个改进方法，使其能适应如用于非线性
所有数据。

采用 SVM，选择 γ 核函数，通过对数据映射到高
维空间，最小化高维特征空间中构造出最优分离超
平面，从而把非线性和非线性数据分离。