

*以 score 检验模型好坏

聚类算法
(第三章)

Supervised Learning

多项式
(第四章)

*如果不考虑内存与时间, 适用大数据集



单变量替换 (Chapter 4)



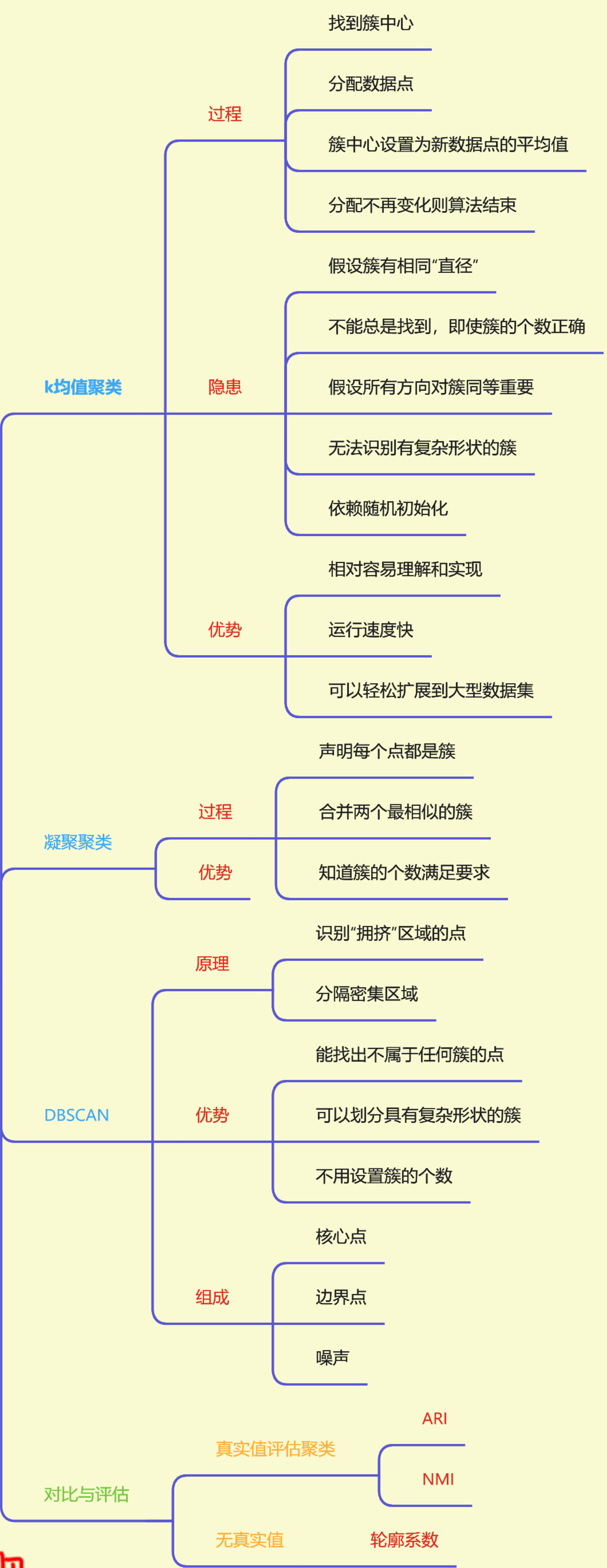
Attention:
缩放器对 train 和 test
作用同样变化

用 fit_transform

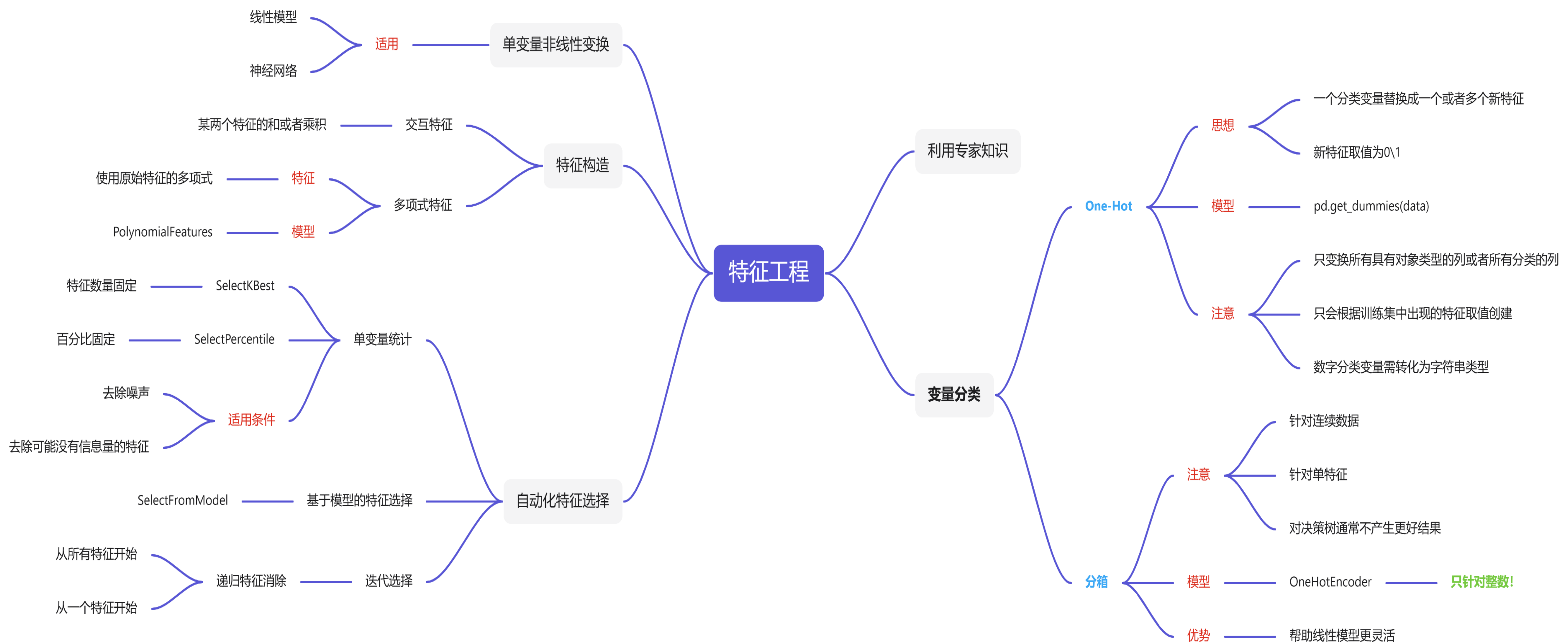
使用 transform

Unsupervised Learning

聚类算法



评价一般
较为主观



- 分箱、多项式, 交互项对复杂度低模型大有提升 线性模型, 朴素贝叶斯模型
- mask: get_support 可视化选中数据

文本数据处理

bag-of-words

导入词袋

只计算每个单词
出现频次

拟合

分词 & 构建词表

CountVectorizer

正则表达式提取词例 $\backslash b \backslash w \backslash w + \backslash b$
所有字母转换为小写 优: 无一个字母的单词
缺: 有很多无信息的特征

Vocabulary-访问词表

transform 训练词袋表示

toarray 查看出现次数

训练分类器

线性模型如
Logistic Regression

优化

min-df

至少在几个文档中出现

stopword
列表

按预计信息量大小
缩放特征

常用

词频-逆向文档频率

在某个特定文档中经常出现则权重高
在多个文档出现则权重低

可使用网格搜索提高准确率

ngram-range=(x, y)

以 x 元分词到 y 元分词

将 replaces, replacement,
replaced, replacing 归为一类

n 元分词

tokenization
optimize

分词优化

词干提取 porter

词形还原 spacy

主题建模

分解

Latent Dirichlet Allocation

查看主题

print-topics

document-topics

求和看到每个主题
整体权重