

第 1 讲	数据科学概述
第 2 讲	数据预处理
	1.数据预处理
	2.数据清洗
	(1)缺失值处理
	(2)噪声平滑
	3.数据集成
	4.数据归约
	5.数据变换
	(1)规范化
	①Min-Max 标准化(0-1 标准化)
	②小数定标标准化
	③Z-Score 标准化
	④Logistic 标准化
	(2)数字编码
	①One-hot 编码
	②TF-IDF
	(3)离散化
	①无监督离散化
第 3 讲	回归与分类模型
第 4 讲	聚类模型
第 5 讲	模型的评价
第 6 讲	集成算法
*第 7 讲	自然语言理解与文本分析
*第 8 讲	社交网络分析

目录

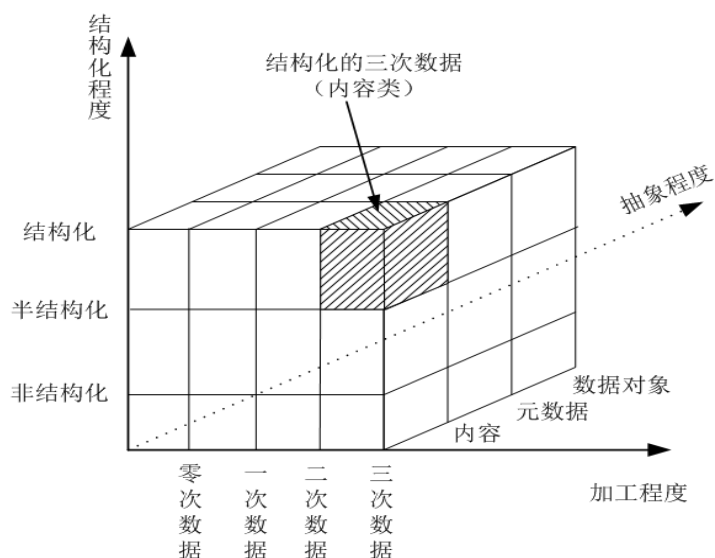
数据科学概述.....	5
1.数据的纬度.....	5
2.数据的结构化程度.....	5
3.数据的加工纬度.....	5
4.大数据.....	6
5.常见的数据类型.....	6
6.数据分析的主要困难.....	6
数据预处理.....	8
1.脏数据.....	8
2.数据集.....	8
3.统计学规律.....	8
4.语言学规律.....	9
5.数据预处理的主要任务.....	9
描述性数据.....	10
1.数据度量的方法.....	10
2.截断均值.....	10
3.中心趋势度量.....	10
4.正态分布函数曲线 $N \sim (\mu, \sigma^2)$	10
缺失值处理.....	11
1.处理方法.....	11
噪声平滑.....	12
1.噪声 Noise	12
2.数据平滑.....	12
3.异常值的检测(识别离群值).....	12
4.异常值的处理.....	14
数据集成.....	16
1.数据冲突.....	16
2.冗余数据检测——相关性分析.....	16
数据归约.....	18
1.三种策略.....	18
2.维归约.....	18
3.数量归约.....	19
4.数据压缩.....	19

数据转换编码.....	20
1.Min-Max 标准化(0-1 标准化)	20
2.小数定标规范化.....	20
3.Z-Score 标准化.....	20
4. Logistic 标准化	21
5.四种方法的对比与一些其他方法.....	21
6.数字编码.....	21
7.One-hot 编码	22
8.离散化.....	23
9.标称数据的概念分层.....	26
先谈机器学习.....	28
一、机器学习.....	28
分类与回归.....	33
0.说在前头:	33
1.分类-决策树	33
2. 分类-k 邻近算法(KNN).....	36
3.分类-朴素贝叶斯	40
4.分类- SVM 支持向量机 Support Vector Machine.....	41
5.回归.....	44
6.关联规则挖掘.....	50
7.推荐系统.....	53
聚类分析.....	54
0.说在前头.....	54
1.划分聚类 (基于原型的聚类 Prototype-Based Clustering).....	54
2.层次聚类.....	56
3.密度聚类.....	57
集成模型.....	59
1.定义与作用与一些术语.....	59
2.如何得到结果.....	59
3.典型的集成方法.....	59
神经网络.....	63
1.神经元模型.....	63
2.感知机.....	63
3.(多层)神经网络 Neural Network	64
4.误差逆传播算法.....	65
5.深度学习	65
6.神经网络一些细节.....	66
7.机器学习&神经网络&深度学习的关系图	66

CNN	67
1.卷积核 convolution kernel/滤波器 filter.....	67
2.多输入多输出.....	67
3.池化/汇聚 pooling	67
RNN	69
PageRank 算法	70
基本思想:	70

数据科学概述

1.数据的纬度



2.数据的结构化程度

类型	含义	本质	举例
结构化数据	直接可以用传统关系数据库存储和管理 ^[1]	先有结构，后有数据	关系型数据库中的数据
非结构化数据	无法用关系数据库存储和管理的数据	没有（或难于发现）统一结构	语音，图像文件
半结构化数据	经过一定的转换 ^[2] 可以用传统关系数据库存储和管理的数据	先有数据，后有结构	HTML,XML 文件 ^[3] 等

【注 1】：传统的关系数据库：如 MySQL、Oracle。

为什么是“直接”：因为结构化数据具有明确定义的模式和规则，可以被分解为表格的形式，以便于使用 SQL 查询语言进行访问和操作。

【注 2】：经过解析和提取，将其中的关键信息转换为结构化数据的形式。

例如，对于 HTML 或 XML 文件，可以使用解析器将其转换为关系型数据库中的表格形式，将不同的标记或元素映射到表的列和行，并根据数据内容创建数据库模式。这种转换过程可以使用各种技术和工具进行，如 XPath、XSLT、正则表达式等。转换后的数据可以进行传统数据库操作，比如查询、插入、更新和删除等。

【注 3】：HTML(Hypertext Markup Language)是用于创建网页的标记语言，描述网页的结构和内容。XML(eXtensible Markup Language)是通用的标记语言，描述文档的结构和语义。

3.数据的加工纬度

零次数据：原始数据（没有经过预处理）

一次数据：加工数据（预处理过的数据）

二次数据：增值数据（分析处理的结果）

三次数据：洞见数据（直接可用于决策）

4. 大数据

大数据依旧是数据，或数据相关的过程，其次，大数据的规模并非一定要达到某一确切的数值，关键在于，是否超过了实际情况下的数据存储能力和数据计算能力。

5. 常见的数据类型

表数据：以行和列的形式组织的结构化数据。行代表**记录**（数据对象），列代表**属性**（对象特征）。比如：数据库表格、电子表格或 CSV 文件。是最为经典的数据。

点集：由一组**坐标或位置信息**组成的数据。比如：地理位置、坐标、传感器数据等。很多数据都可以看成是某种空间的点的集合。

时间序列：按**时间顺序**排列的数据集合，每个数据点都与特定时间点/时间段相关联。例如：温度数据、文本、通话、DNA 序列。常用于分析和预测趋势、周期性和季节性等。

图像视频：由像素组成的数据集。可以看成两个变量(空间位置和时间)的函数。

网页和报纸：非结构化或半结构化的数据。每篇文章都可以看成是时间序列，整个网页或报纸又具有空间结构。

网络数据：描述网络结构和连接关系的数据。网络数据本质上是图，由节点和联系节点的边构成，节点表示实体（如电脑、服务器、人物等），边表示节点之间的关系或连接（如通信链路、社交关系等）。

6. 数据分析的主要困难

(1) **数据量大**

(2) **维数高**：模型复杂度和计算量随着维数的增加而指数增加。
将模型限制在一个极小的特殊类里面，如线性模型。
利用数据可能有的特殊结构，例如稀疏性，低秩，光滑性等，通过各类降维方法。

(3) **类型复杂**：结构化、半结构化、非结构化

7. 算法处理大数据的两条思路

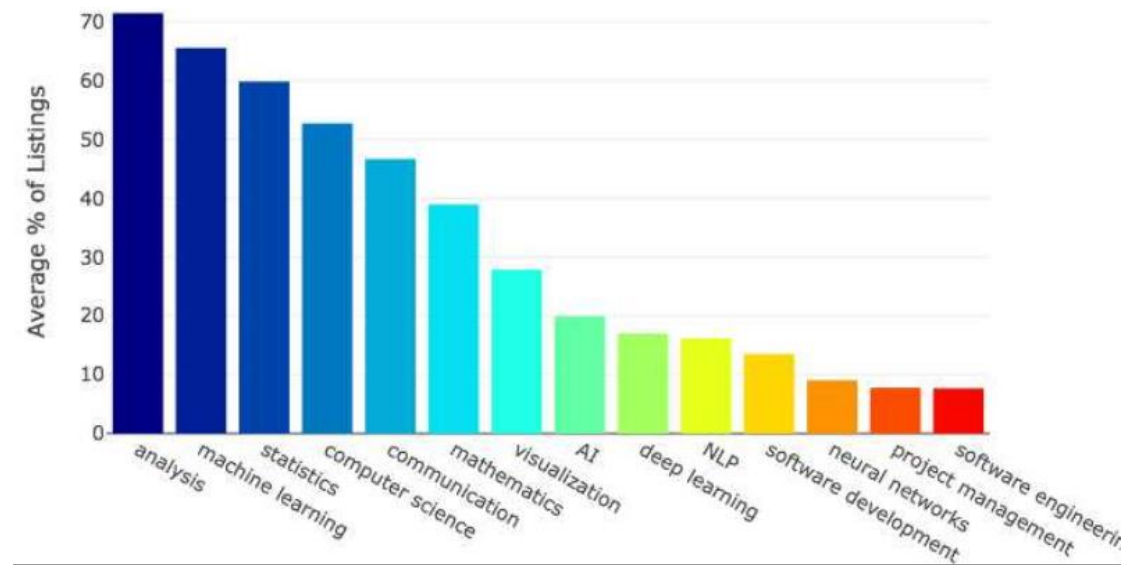
(1) 降低算法的**复杂度**；

(2) 分布式计算，**分而治之**，如著名的 MapReduce 框架

8. 常用工具

R、Python、Clojure、Haskell、Scala 等数据科学语言工具；
NoSQL、MongoDB、Couchbase、Cassandra 等 NoSQL 工具
SQL、RDMS、DW、OLAP 等传统数据库和数据仓库工具；
HadoopHDFS&MapReduce、Spark、Storm 等支持大数据计算的工具；
HBase、Pig、Hive、Impala、Cascalog 等支持大数据管理、存储和查询的工具。
Webscraper、Flume Avro、Sqoop、Hume 等支持数据采集、聚合或传递的工具
Weka、Knime、RapidMiner、SciPy、Pandas 等支持数据挖掘的工具
ggplot2、D3.js、Tableau、Shiny、Flare、Gephi 等支持数据可视化的工具
SAS、SPSS、Matlab 等数据统计分析工具

9.general skill in data scientist job listings



数据预处理

1.脏数据

不完整: 缺少数据值; 缺乏某些重要属性; 仅包含汇总数据(原因:并未填写)

有噪声: 包含错误或者孤立点(原因:故意输错;不愿提交)

不一致: 在编码或者命名上存在差异, 如: 等级: 1、2、3; A、B、C; 甲、乙、丙
重复记录间的不一致性(函数依赖性), 如: Age="20" Birthday="03/07/1997"
(原因:数据传输错误;命名规范/格式不一致)

2.数据集

数据集由**数据对象**组成。一个数据对象代表一个实体, 由**属性**来描述。

数据库的**行对应数据对象**, **列对应属性**。

属性

表示数据对象的特征。

在文献中, **属性**、**维度** (dimension)、**特征** (feature)、**变量** (variance) 可以互换的使用。

“维”: 一般用在数据仓库中; “特征”: 一般用在机器学习中; “变量”: 一般用在统计学中。

标称属性: 可枚举的(分类的)。用数来这些符号或名称, 所以均值与中位数没有意义(但众数有意义)。为避免编码时产生不同的距离, 可以这样编码:A(0,0,1) B(0,1,0) C(1,0,0)

二元属性: 0/1, 也称布尔属性。如性别。

序数属性: 值是有意义的序, 但相继值之间的差是未定义的。如大小、胜负、优良差。
可以用众数中位数表示序数属性的中性趋势, 但均值没有意义。

数值属性: 不同于上面三个定性的量, 数值属性是可度量的量, 比如销量。

离散属性: 有限个或者可数无限个, 如 ID、颜色

连续属性: 不是离散的属性, 如成绩。与“数值属性”可以互换使用。

3.统计学规律

第一数字定律 Benford's Law:

在 b 进位制中, 以数 N 起头的数出现的概率为 $P(N) = \log_b(N+1) - \log_b(N) = \log_b \frac{N+1}{N}$ 。

数字 N	1	2	3	4	5	6	7	8	9
以 N 为首位的概率%	30.1	17.6	12.5	9.7	7.9	6.7	5.8	5.1	4.6

适用条件: 数据不能经过人为修饰; 数据不能是规律排序, 比如发票编号, 身份证号码。

用途: 可以用于判断数据是否**可能有问题**。

小概率原理 Law of Small Probabilities:

一个事件如果发生的概率很小的话, 那么它在一次试验中是几乎不可能发生的, 但在多次重复试验中几乎是必然发生的, 数学上称之小概率原理。在统计学中, 把小概率事件在一次实验中看成是实际不可能发生的事件, 一般认为等于或小于 0.05 或 0.01 的概率为小概率。

齐夫定律 Zipf's law:

在自然语言的语料库里, 一个单词出现的频率与它在频率表里的排名成反比。是幂律分布的一种特殊形式。

用途: 单词的出现频率、网页访问频率、收入前 3% 的人的收入、地震震级、固体破碎

时的碎片大小。

4.语言学规律

连接性特征:

包括语言学中的后连接(如字母“q”后总是“u”)、前连接(如字母“x”的前面总是字母“i”,字母“e”很少与“o”和“a”连接)以及间断连接(如在“e”和“e”之间,“r”的出现频率最高)。

重复特征:

两个字符以上的字符串重复出现的现象,叫做语言的重复特征。例如在英文中字符串“th”、“tion”和“tious”的重复率很高。

5.数据预处理的主要任务

(1)描述性数据归总

理解数据的基本情况,发现异常值、检查数据可靠性和一致性

数据清洗: 填写空缺的值,平滑噪声数据,识别、删除孤立点,解决不一致性

数据集成: 集成多个数据库、数据立方体或文件

(2)数据变换和标准化

规范化和聚集,提高涉及距离度量的挖掘算法的准确性和有效性

数据归约: 得到数据集的压缩表示,它小得多,但可以得到相同或相近的结果

数据离散化: 数据归约的一部分,通过概念分层和数据的离散化来规约数据,对数字型数据特别重要

(3)特征的选择与降维

选择最相关的特征,以提高模型效率,减少数据的复杂性

①**特征选择:** 从所有特征中选择最具有预测能力的特征子集,以提高模型的准确性、降低过拟合的风险、减少计算成本。

过滤法: 通过计算特征和目标变量之间的相关性或统计指标来进行特征选择。

包裹法: 使用模型训练和评估的结果来进行特征选择。

嵌入法: 将特征选择嵌入到机器学习算法的训练过程中,通过正则化或特征权重更新来选择重要的特征。

②**降维:** 减少数据的维度,以消除冗余和噪声,并简化模型的复杂性。

主成分分析 PCA: 通过线性变换将原始特征映射到新的正交特征空间,以捕捉数据中的主要方差。

线性判别分析 LDA: 是一种有监督的降维方法,它通过最大化类别之间的方差和最小化类别内部的方差来找到具有最佳分类能力的特征投影。

特征抽取: 一种非线性的降维方法,它利用神经网络或自编码器等技术将原始特征映射到低维空间。

描述性数据

1. 数据度量的方法

度量数据的**中心趋势**：均值、中位数、众数、中列数

度量数据的**离散程度**：四分位数、四分位数极差、方差

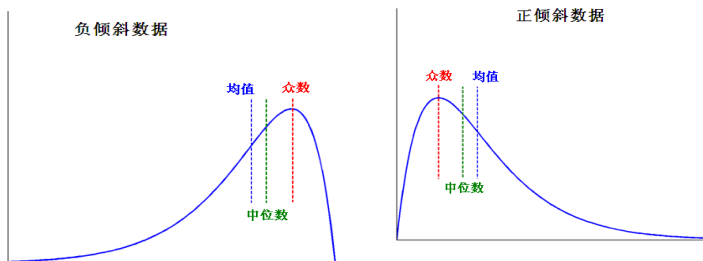
2. 截断均值

去掉高、低极端值得到的均值。

如：计算平均工资时，可以截掉上下各 2% 的值后计算均值，以抵消少数极端值的影响。

3. 中心趋势度量

左右对称：均值=中位数=众数 **负倾斜**：均值<中位数<众数 **正倾斜**：众数<中位数<均值



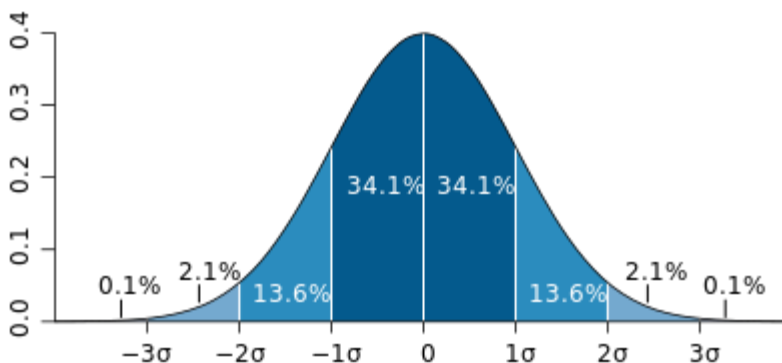
众数与中位数的距离约为中位数与均值距离的 2 倍。

某个异常值对**均值**的影响最大(均值敏感)

4. 正态分布函数曲线 $N \sim (\mu, \sigma^2)$

正态分布 Normal distribution 也叫高斯分布 Gaussian distribution

$(\mu - \sigma, \mu + \sigma)$: 86.27% , $(\mu - 2\sigma, \mu + 2\sigma)$: 95.45% , $(\mu - 3\sigma, \mu + 3\sigma)$: 99.73%



缺失值处理

1.处理方法

(1)删除带有缺失值的样本或特征:

删除样本: 删除存在数据缺失的样本。该方法适合某些样本有多个特征存在缺失值,且存在缺失值的样本占整个数据集样本数量的比例不高时的情形。

删除特征: 当某个特征缺失值较多,且该特征对数据分析的目标影响不大时,可以将该特征删除。

(2)采用某种方法对缺失值进行填补:

均值填补: 连续型特征,平均值填补;离散型特征,众数填补。但会降低方差、相关性。如果数据是倾斜的,应该使用中位数,或者同一类样本的中位数。

适用于对个体精度不大的数据。

实际可以根据一定的辅助特征,先将数据集分组,然后在每一组数据上分别使用均值填补,比如学生信息数据集,"入学年份"特征记录了学生入学的年份,可以首先根据毕业年份对数据进行分组,然后使用每一个分组数据中"年龄"特征的平均值来对缺失值进行填补。

随机填补: 随机填补是在均值填补的基础上加上随机项,通过增加缺失值的随机性来改善缺失值分布过于集中的缺陷。

贝叶斯 Bootstrap 方法: 设数据集有 n 个样本,某特征 f 存在 $n-k$ 个缺失值。

① 先从均匀分布 $U \sim (0,1)$ 中随机抽取 $k-1$ 个随机数并升序排序为 $\{0, a_1, \dots, a_{k-1}, 1\}$ 。

② 对 $n-k$ 个缺失值,分别从非缺失值 $\{f_1, f_2, \dots, f_k\}$ 中以概率 $a_1-0, a_2-a_1, \dots, 1-a_{k-1}$ 采样一个值进行填补。

近似贝叶斯方法:

① 从非缺失值 $\{f_1, f_2, \dots, f_k\}$ 中有放回地抽取 k 个值建立一个新集合 F 。

② 对 $n-k$ 个缺失值,分别从 F 中随机抽取一个值进行填补。

基于模型的填补: 由非缺失值构建训练集,训练模型来预测缺失值。比如根据父母身高推测丢失的子女身高。实际使用时需要采用模型评估方法对模型的预测性能进行评估,如果构建的模型预测性能太差,则不适合使用该方法。此外,基于模型的填补方法将增大特征之间的相关性。

其他填补法:

全局常数(哑变量方法):

将缺失的属性值用同一个常数(如 Unknown 或 NaN)替换。但比如将性别特征的缺失值作为一个特殊的取值"unknown",表示性别未知。此时认为"性别"特征包含"F"、"M"和"unknown"三个不同取值。

EM 算法:

利用不完整的信息实现概率模型的参数化估计的算法。该算法可以用来进行缺失值填补,此时缺失特征被当作隐含变量。

2.文献阅读:

(1)缺失数据, 保罗·D.埃里森 林毓玲译出版社:格致出版社出版, 2018.06

(2)基于缺失值变量的不完整数据填补与分类研究, 硕士学位论文, 大连理工大学, 吴霞, 2020.

噪声平滑

1. 噪声 Noise

噪声是被测量的变量的随机误差或方差，包括错误的值或偏离期望的离群点。

2. 数据平滑

分箱法：

思想：考查数据的近邻来平滑有序数据值。是**局部平滑**。

方法：将数据划分成几个含有相同个数数据的数据段，再将箱中的**每个值用某个值替换**（减少值的类别数量）。一般而言，箱的宽度**越大**，平滑效果**越明显**。

某个值的选择：均值、中位数、边界（将箱中的最大/小值试作边界，每个值被最近的边界值替换）。

方法与理论知识：

移动平均 Moving Average：计算数据点周围一定窗口大小内的均值来减少噪声的影响。

中值滤波 Median Filtering：非线性噪声平滑，使用数据点周围窗口内的中值来替代当前数据点的值。常用于去除椒盐噪声(图像中的随机黑点和白点)等脉冲噪声。

加权平均 Weighted Averaging：分配不同权重，通常在窗口中心的数据点权重较高，而边缘的权重较低。

指数平滑 Exponential Smoothing：较新的数据点权重较高，较旧的权值较低。常用于时间序列数据(用于捕捉趋势和季节性变化)。

例子：设有数据，升序排列后为 4,8,15,21,21,24,25,28,34，将其划分为三个箱子。

用箱均值平滑：箱 1：9,9,9。箱 2：22,22,22。箱 3：29,29,29。

用箱边界平滑：箱 1：4,4,15。箱 2：21,21,24。箱 3：25,25,34。

用箱中值平滑：箱 1：8,8,8。箱 2：21,21,21。箱 3：28,28,28。

回归法：

用一个函数拟合数据来平滑数据。

离群点分析：

通过聚类来检测离群点，将落在簇集合之外的值视为离群点。

3. 异常值的检测(识别离群值)

简单统计：

判断变量的取值是否在合理的范围（如年龄，负数或者 200 岁肯定是异常值）。

3 σ 原则：

参见[正态分布](#)。

箱型图 Tukey's test:

百分位数 percentile：第 k 个百分位数是具有如下性质的值 x：k% 的数据项位于或低于 x，中位数就是第 50 个百分位数

四分位数：

Q1 (下四分位数，也称 QL，25th percentile，第 $(n+3)/4$ 个数)

Q3 (上四分位数，也称 QU，75th percentile，第 $(3n+1)/4$ 个数)

中间四分位数极差/四分位间距(interquartile range, IQR): $IQR = Q3 - Q1$

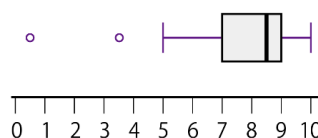
孤立点/离群值(outlier): 落在至少高于第三个四分位数或低于第一个四分位数 $1.5 \times IQR$ 处的值。位于范围外 $1.5 \times IQR$ 到 $3 \times IQR$ 范围的数值, 称作适度离群值 (mild outlier)。位于范围外 $3 \times IQR$ 以上的数值, 称作极端离群值 (extreme outlier)。

上下限: 上限= $\min\{Q3+1.5IQR, \text{MAX}\}$, 下限= $\max\{Q1-1.5IQR, \text{MIN}\}$ 。

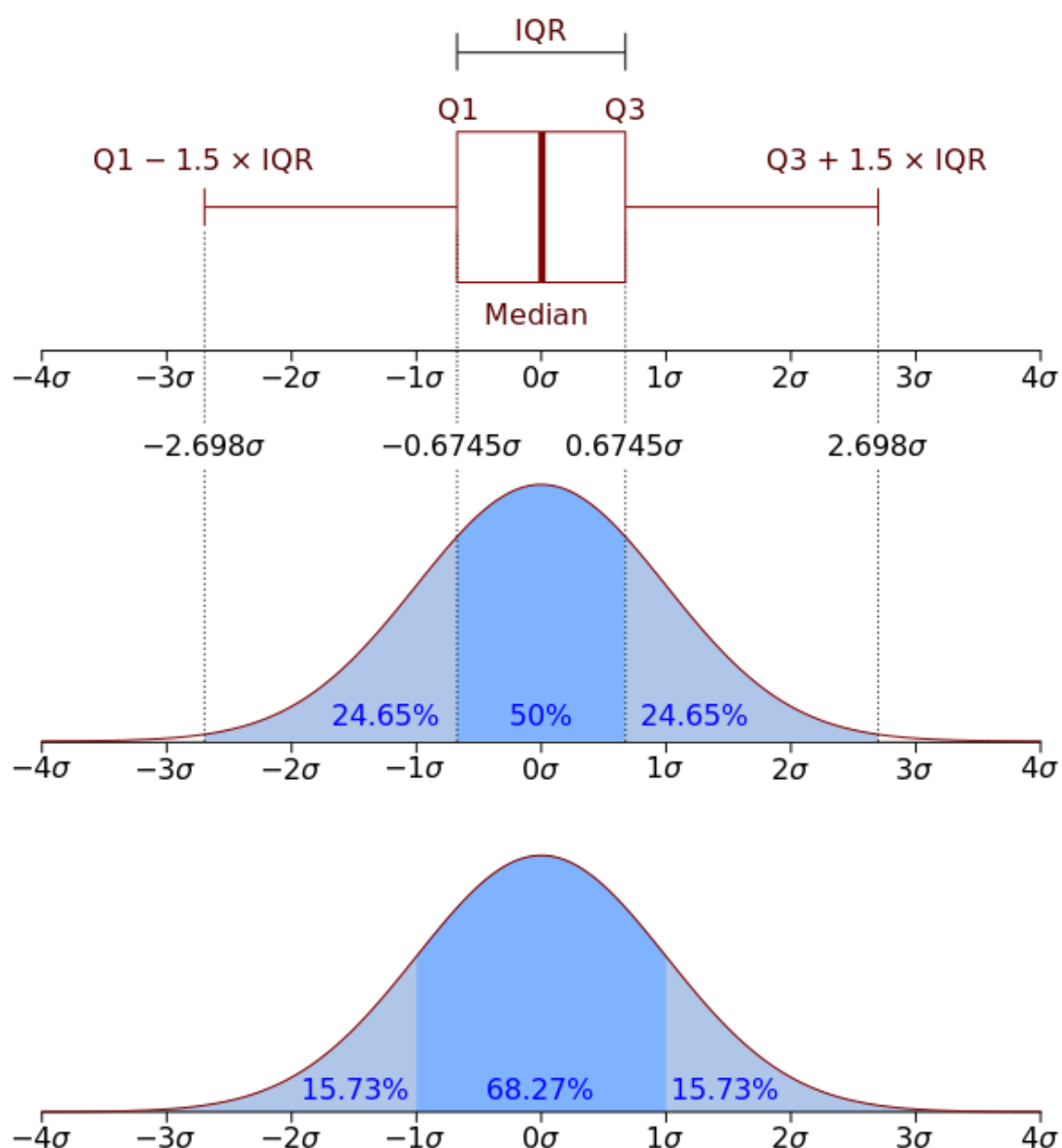
换句话说, 上下限是忽略掉异常值 (包括温和异常值) 之后的最大值和最小值。

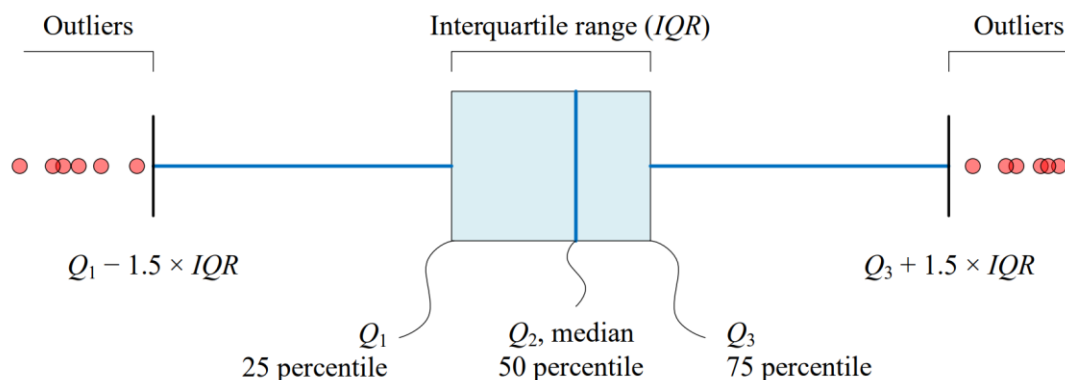
如右图, 下边界=5, $Q1=7$, $Q2=8.5$, $Q3=9$, 上边界=10,

$IQR=Q3-Q1=2$ (即 ΔQ)。适度离群值=3.5, 极端离群值=0.5。



箱型图与正态分布:





基于模型检测:

如果是建立数据模型, 异常值是不能完美拟合模型的对象。

如果模型是簇的集合, 异常值是不显著属于任何簇的对象。

如果是使用回归模型, 异常值是相对远离预测值的对象。

优势: 数据充分、检验知识完备的情况下。

劣势: 对于多元数据、高维数据的性能不好。

基于距离:

定义邻近性度量, 异常对象就是那些远离其他对象的对象。

优点: 简单。

缺点: 时间复杂度 $O(n^2)$ 且 对参数的选择敏感 且 不能处理不同密度区域的数据集(因

为它使用的是全局阈值, 不能考虑密度的变化)

基于密度:

当一个点的局部密度显著低于它的大部分近邻时才将其分类为离群点。

优点: 适合非均匀分布的数据。给出了离群点的定量度量。即使数据具有不同的区域也能很好处理。

缺点: 时间复杂度 $O(n^2)$ 且 参数选择困难。

基于聚类:

当对象不强属于任何簇时, 它是基于聚类的离群点。

优点: 线性复杂度。可同时发现离群点和簇。

缺点: 非常依赖所用的簇的个数和数据中离群点的存在性。簇的质量对离群点的质量影响非常大。

4.异常值的处理

- 1.删除异常值
- 2.视为缺失值处理
- 3.用均值修正
- 4.不处理

数据集成

1. 数据冲突

背景:

信息系统独立。数据源分布异构等原因形成了“信息孤岛”。

数据集成就是合并来自多个数据源的数据。

产生原因:

(1)同名异义:比如数据源 A 中的属性 ID 与 B 中的 ID 不是一个 ID(描述不同的实体)。

(2)异名同义:比如数据源 A 中的 sales_dt 与 B 中的 sales_data 是一样的。

解决方法:

利用元数据 metadata (data about data) 来解决。

每个属性的元数据包括名字、含义、数据类型、取值范围、处理空值的规则、其他描述信息等。以学生信息数据集为例:

属性可以是学生的姓名、年龄和性别。

元数据则可以包括属性的定义、数据类型(字符串、整数等)、取值范围(如年龄 1~100)、单位(如年龄用岁表示)以及其他描述信息(如年龄是指学生的当前年龄,而非出生年月日)。

也就是说,属性是数据实际具有的特征或特性,而元数据是描述这些属性的数据。

2. 冗余数据检测——相关性分析

冗余数据:属性 B 可以由属性 A 导出,则属性 B 可能是冗余的。(比如出生导出年龄)

对于标称数据,可以用 χ^2 卡方检验。

对于数值属性,可以用相关系数 Correlation Coefficient 和协方差 Covariance。

χ^2 相关性检测(卡方检验):

设 A 有 c 个不同值 a_1, a_2, \dots, a_c 构成列, B 有 r 个不同值 b_1, b_2, \dots, b_r 构成行。记 A 的第 i 个值与 B 的第 j 个值构成的联合事件为 (A_i, B_j) 。

$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(A_{ij} - T_{ij})^2}{T_{ij}}$, 式中 A_{ij} 是联合事件 (A_i, B_j) 的观测频度,即实际计数 actual, T_{ij} 是

(A_i, B_j) 的期望频度 Theoretical, 即 $(A$ 中 a_i 的个数 \times B 中 b_j 的个数) $\div n$ 。然后依据自由度与显著性水平查表确定拒绝域的临界值(自由度 $= (r-1) \times (c-1)$, 显著性水平常取 $\alpha=0.05$)

例如右图, a, b, c, d 为观测频度, 括号内为期望频度。

$$\chi^2 = \sum_{i=1}^2 \sum_{j=1}^2 \frac{(A_{ij} - T_{ij})^2}{T_{ij}} \sim \chi^2(2)$$

性别\视觉	正常	色盲	总和
男	$a \left(\frac{(a+c)(a+b)}{n} \right)$	$b \left(\frac{(b+d)(a+b)}{n} \right)$	$a + b$
女	$c \left(\frac{(a+c)(c+d)}{n} \right)$	$d \left(\frac{(b+d)(c+d)}{n} \right)$	$c + d$
总和	$a + c$	$b + d$	n

协方差(Covariance)与相关系数(Pearson 积矩系数):

协方差 $\text{Cov}(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$, 其中 \bar{A}, \bar{B} 相当于 EA, EB , 公式

简化后为 $\text{Cov}(A, B) = E(AB) - E(A)E(B)$ 。

相关系数 $\rho_{A,B} = \frac{\text{Cov}(A, B)}{\sigma_A \sigma_B}$ 。 $\rho_{A,B} = 0$ 表示 A, B 不相关(但可能存在线性关系以外的其他

函数关系), $\rho_{A,B} > 0$ 表示 A, B 正相关。

相关性，描述两个随机变量是否存在线性关系。**独立性**，则考察的是两个随机变量是否存在某种关系。因此，**相互独立事件必然不相关，反过来却不一定成立**。此外，若随机变量 X, Y 服从二元高斯分布，则二者不相关与独立互为充要条件。相关性不蕴含因果性。

三种相关系数的公式计算可参考 <https://zhuanlan.zhihu.com/p/431865174>。

协方差的理解可参考 <https://www.zhihu.com/question/20852004> 中 GRAYLAMB 的回答。

元组重复：

对于给定的唯一数据实体，存在两个或多个相同的元组。

去规范化表(denormalized table)的使用（这样做通常是通过避免连接来改善性能，创建的是宽表）是数据冗余的另一个来源。不一致通常出现在各种不同的副本之间，由于不正确的输入，或者由于更新了部分相关数据，但是未更新所有的相关数据。

小结：

数据冗余的可能情况：

- 1.同一属性多次出现或仅是命名前后不一致。
- 2.A 属性能由 B 导出。

数据归约

1. 三种策略

数据归约：用数据集的规约表示，通过比原数据集小很多的数据产生几乎相同的分析结果。有三种方法：

维归约 Dimensionality Reduction：减少数据集的**特征维度**。将原始数据中的冗余或不重要的特征删除，同时保留对数据中重要信息的捕获。

数量归约 Numerosity Reduction：减少数据集中的行数或**样本数量**。在不损失太多信息的情况下，减少数据集的规模，以加快分析和处理的速度。

数据压缩 Data Compression：使用**压缩算法**来减小数据集的大小，同时尽量保留数据的信息。压缩方法的思想是用更紧凑的表示方式来存储数据，从而减小存储空间。

2. 维归约

主成分分析 Principal Component Analysis(PCA)：又称 K-L 方法，搜索 k 个最能代表数据的 n 维正交向量($k \leq n$)，使得数据投影到一个小得多的空间上(新的坐标系)，实现维规约。

过程：

①对输入数据 $X = \{x_1, x_2, \dots, x_n\}$ 规范化 $x_i' = x_i - \frac{\sum_{j=1}^n x_j}{n}$ ，使得每个属性都落入相同的区间(避免较大定义域的属性不会支配较小定义域的属性)

②计算协方差矩阵 $\frac{XX^T}{n}$ ，求特征值与特征向量，选择最大的 k 个特征值(作为规范化输入数据的基)。然后将这 k 个特征值对应的 k 个特征向量 w_{1-k} 标准化后分别作为列向量组成特征向量矩阵 W 。

③将数据映射到 k 个特征向量构建的新空间中 $z_i = W^T x_i$ ，输出数据 $X' = (z_1, z_2, \dots, z_n)$ 。

主成分到底怎么理解：

主成分其实就是数据的**新坐标系**，提供关于方差的信息。第一个轴显示的数据方差最大，第二个次之，如此下去。

适用情景：

有序、无序的属性，稀疏、倾斜数据等。将多维数据归约为二维数据。作为多元分析和聚类分析的输入。

特征选择 Feature Selection：选择原始数据集中的部分特征而不创建新的特征空间。常用的方法包括基于统计测试的方法、递归特征消除 RFE。

选择方法：

暴力选择需要 $O(2^n)$ ，这里有四种方法：

逐步向前选择：从空属性集开始，每次选择当前最好的属性(使用统计显著性测试)。

逐步向后删除：从全集开始，每次删除当前最坏的属性。

逐步向前选择和逐步向后删除结合：每一步选择一个最好的属性并删除一个最坏的属性。

决策树归纳：使用决策树生成的内部节点表示属性(树叶对应类预测)。

线性判别降维法 Linear Discriminant Analysis(LDA)：找到在不同类别之间具有最大可分性的投影方向。不同于 PCA 方差最大化理论，LDA 算法希望将数据投影到低维空间之

后，同一类数据尽可能的紧凑，不同类的数据尽可能分散。

3.数量归约

随机抽样：从数据集中随机选择样本，以减少数据的大小。还可以分层抽样。

聚类采样：将相似的数据点分组在一起，然后通过单个代表性数据点表示每个组

过滤异常值：移除数据中的异常值，以减小数据集的噪音和离群点。

4.数据压缩

有损压缩：有损压缩方法会删除一些数据，以减小文件大小，但可能损失一部分信息。这种方法常用于图像、音频和视频压缩，例如 JPEG 图像压缩。

无损压缩：无损压缩方法可以压缩数据，但在解压后不会损失任何信息。这种方法用于需要完整保存数据的场合，例如 ZIP 压缩。

数据转换编码

1.Min-Max 标准化(0-1 标准化)

通过线性变换 $x_i^* = \frac{x_i - \min}{\max - \min}$ 使处理过后数据均落在 $[0,1]$ 区间内。

通过线性变换 $x_i^* = \frac{x_i - \frac{\max + \min}{2}}{\frac{\max - \min}{2}}$ 使处理过后数据均落在 $[-1,1]$ 区间内。

适用范围：

0-1 标准化适用于需要将数据简单地变换映射到某一区间中，但其不足之处在于当有新数据加入时，可能会导致数据系列中的最大值或最小值发生变化，此时便需要重新定义最大值、最小值(也就是需要重新标准化)，否则就会越界。

如何把某数据集映射到指定的区间 $[a,b]$ ： $x_i^* = (b-a) \frac{x_i - \min}{\max - \min} + a$

2.小数定标规范化

通过移动数据的小数点位置来进行标准化，小数点移动多少位取决于数据系列中的最大

绝对值大小： $x_i^* = \frac{x_i}{10^j}$ (j 是使得 $\max(|x_i^*|) < 1$ 的最小整数)

适用于数据系列分布比较离散，尤其是数据系列遍布多个数量级的情况，若数据系列分布集中在某几个量级上，则小数定标标准化的数据也会集中在某个值附近，不利于建模时的数据区分。优点在于简单实用，在确定了小数点移动位数后，易于还原标准化后的数据系列。

3.Z-Score 标准化

最常用的标准化方法，使得处理后的数据具有固定均值0和标准差1(但不一定是正态)。

基于平均值与标准差规范化： $x_i^* = \frac{x_i - \mu}{\sigma}$ 。其中 $\mu = \frac{\sum_{i=1}^n x_i}{n}$, $\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{n}}$

当数据中存在离群点时，为了降低离群值的影响，将 σ 替换成平均绝对差 $s = \frac{\sum_{i=1}^n |f_i - \mu|}{n}$

$|x_{Z-Score}| > 3$ 的样本认为是离群值。

经过 Z-Score 标准化后的数据，能直观反映每个数据点距离平均值点的标准差距离，从而理解整体数据的分布情况。每个数据点离零点的距离可解释为其远离均值的标准差距离。

适用范围：数据系列中最大值或最小值未知、数据分布非常离散(离群点左右了 Min-Max 标准化)的情况。

应用场景：通过将数据标准化到同一比例，可以更好地比较不同变量之间的关系，或者比较来自不同分布的数据。设 A 的满分是 100 分，B 的满分是 700 分，显然 A 考 70 分与 B 考 70 分代表的意义完全不同。因此需要 Z-Score(一个同等的标准)来比较 A 与 B 的成绩。

4. Logistic 标准化

利用 Sigmoid 函数(也称 Logistic 函数) $x_i^* = \frac{1}{1 + e^{-x_i}}$, 将原始数据系列转化为 [0,1] 之间。

适用范围: Logistic 标准化方法适用于数据系列分布相对比较集中地分布于零点两侧, 否则容易被压缩到两端。

5. 四种方法的对比与一些其他方法

对比表格:

方法	优点	缺点	适用范围
0-1	对数据作线性变换, 保留数据的原始关系	若原数据最大最小值发生变化需重新定义	需保留原始数据间的关系, 且最大/小值已经确定
小数定标	简单实用、易于还原标准化后的数据	若原数据最大绝对值发生变化需重新定义	数据系列分布比较离散, 数据遍布多个数量级
Z-Score	转化为标准正态分布无序数据的最大、最小值	需要记录原数据均值方差	数据中最大最小值未知且数据系列分布离散
Logistic	简单, 通过单一映射函数对数据进行标准化	对分布零散且远离零点的数据处理效果不佳	数据系列分布比较集中, 且均分布于零点两侧

其他方法:

1. 均值归一法 Mean normalization:

$$\textcircled{1} x_i^* = \frac{x_i - \mu}{\max} \quad \textcircled{2} x_i^* = \frac{x_i - \mu}{\max - \min} \quad \text{式中 } \mu = \frac{\sum_{i=1}^n x_i}{n}$$

2. 向量归一化:

$$x_i^* = \frac{x_i}{\sum_{i=1}^n x_i}$$

3. 对数/指数转换:

除了 Logistic 外还有:

$$\textcircled{1} \lg \text{函数: } x_i^* = \frac{\lg x_i}{\lg \max}, \text{ 要求数据都大于 1, 结果就会落入 } [0,1] \text{ 区间内。}$$

$$\text{或: } x_i^* = \log(x_i + 1)$$

$$\textcircled{2} \text{softmax 函数: } x_i^* = \frac{e^{x_i}}{\sum_{i=1}^n e^{x_i}}$$

6. 数字编码

数据源: UCI 的 Wine Data Set。 <http://archive.ics.uci.edu/ml/datasets/Wine>

与 <http://archive.ics.uci.edu/ml/datasets/car+evaluation>

如: “汽车品牌”={路虎, 吉利, 奥迪, 大众, 奔驰}, 数字编码后转换成“汽车品牌”={0, 1, 2, 3, 4}。“收入水平”={贫困, 低收入, 小康, 中等收入, 富有} 数字编码后转换成“收入水平”={0, 1, 2, 3, 4}。

存在的缺陷: 会引入错误的序; 会引入错误距离

7.One-hot 编码

又称一位有效编码/独热编码，其方法是使用 N 位状态寄存器来对 N 个状态进行编码，每个状态都有他独立的寄存器位，并且在任意时候，其中只有一位有效(为 1)，其余全为 0。可以这样理解，对于每一个特征，如果它有 m 个可能值，那么经过独热编码后，就变成了 m 个二元特征，并且这些特征互斥，每次只有一个激活，因此，数据会变成稀疏的。

如：美国(0,0,1)英国(0,1,0)法国(1,0,0)

优点：

避免不必要的相关性：不会给名义型特征的取值人为地引入次序关系。

避免数值影响：不同的原始特征取值之间拥有相同的距离。

便于处理缺失值：如果某个域缺失了，独热编码会分配一个全零向量。

更好地找到特征的关系：独热编码进行交叉乘积可以明显地得到特征之间的关系。

缺点：

特征维度会显著增多：假设存在 10 个包含 100 个取值的离散型特征，经过 One-Hot 编码后的特征数量将变成 1000 个。

增加特征之间的相关性：如美英法编码后，三个特征存在线性关系： $f_1 + f_2 + f_3 = 1$ 。(所以应

该做出改变：对于一个包含 K 个取值的离散型特征，将其转换成 $K-1$ 个二元特征，即哑变量编码，如

原始特征取值	f_1	f_2	f_3	f_4	f_5
路虎	1	0	0	0	0
吉利	0	1	0	0	0
奥迪	0	0	1	0	0
大众	0	0	0	1	0
奔驰	0	0	0	0	1

原始特征取值	f_1	f_2	f_3	f_4
路虎	1	0	0	0
吉利	0	1	0	0
奥迪	0	0	1	0
大众	0	0	0	1
奔驰	0	0	0	0

One-hot:

哑变量:

一个简单的例子如下，编码前后如图：

	Feature_1	Feature_2	Feature_3
Sample_1	1	4	3
Sample_2	2	3	2
Sample_3	1	2	2
Sample_4	2	1	1

	Feature_1	Feature_2	Feature_3
Sample_1	0 1	1 0 0 0	1 0 0
Sample_2	1 0	0 1 0 0	0 1 0
Sample_3	0 1	0 0 1 0	0 1 0
Sample_4	1 0	0 0 0 1	0 0 1

特征向量就为：

Sample_1 = [0,1,1,0,0,0,1,0,0], Sample_2 = [1,0,0,1,0,0,0,1,0], Sample 3,4 同理。

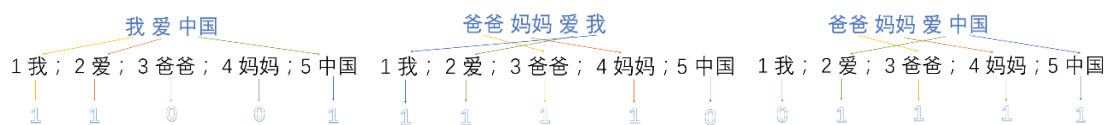
这有什么用？一是解决了分类器不好处理离散数据的问题，二是在一定程度上也起到了扩充特征的作用（上面样本特征数从 3[3 个 Feature]扩展到了 9[9 个向量特征]）。编码是为了分类，虽然这些 Feather 是数字，看似有大小关系，其实只是一种类别，因此可以编码为独立的特征向量。

One-hot 在文本特征提取上属于词袋模型 bag of words，以以下语料库为例：

我爱中国 爸爸妈妈爱我 爸爸妈妈爱中国

编码：1 我 2 爱 3 爸爸 4 妈妈 5 中国

提取特征向量：



因此我们得到了最终的特征向量为

我爱中国 = 1,1,0,0,1 爸爸妈妈爱我 = 1,1,1,1,0 爸爸妈妈爱中国 = 0,1,1,1,1

One-hot 的用处：**特征提取**。

另一种特征提取方法为 **TF-IDF**：

IF-IDF 是信息检索(IR)中最常用的一种文本表示法。统计每个词出现的词频 TF (Term Frequency)，然后乘上权值参数 IDF (逆文档频率 Inverse Document Frequency)。举个例子：

比如统计一篇文档中的前 10 个关键词。首先想到词频越高，这个词就越重要。但关键词最终基本都是“的”、“是”这样没有实际意义的词(停用词)，因此应该加权值：

TF = 某个词出现次数/文章总词数，IDF = $\log[\text{语料库的文档总数}/(\text{包含该词的文档数}+1)]$

此算法无法处理一词多义与一义多词的情况。

8.离散化

算法对数据特征类型通常是有要求的：

关联规则算法只能处理布尔类型的数据，决策树算法只能处理特征为离散型的数据。很多算法能处理连续性特征，但是离散化处理能使得部分模型产生更好的预测效果。

将连续性特征转换为离散型特征的过程称为**特征离散化** data discretization。

将连续性特征的取值范围划分为若干区间段 bin,用区间段代替落在该区间段的特征取值，区间段之间的分割点称为切分点 cut point，分割出来的区间段的个数称为元数 arity。

如何在数据信息损失尽量少的前提下，尽可能减少元数，是数据离散化追求的目标，所以**如何选择切分点**，产生合理的子区间段成为决定特征离散化成败的**关键**。

离散化步骤：

- ①**特征排序**：对连续型特征的取值进行升序或者降序排列，减少离散化的运算开销。
- ②**切分点选择**：根据给定评价准则(比如基于信息增益或者基于统计量)选择切分点。
- ③**区间段分割或合并**：基于选择好的切分点，对现有区间段进行分割或合并，得到新的区间段。
- ④在生成的新的区间段上重复以上步骤，直到满足终止条件。

离散化方法：

自顶向下的离散化方法：等距离散化、等频离散化、基于信息增益的离散化方法。

自底向上的离散化方法：基于卡方的离散化方法。

无监督的离散化方法：等距离散化、等频离散化、基于聚类分析的离散化方法。

有监督的离散化方法：基于信息增益的离散化方法、基于卡方的离散化方法。

等距离散化：

将属性的值域[Min,Max]划分成 K 个区间，每个区间长度相等(均为 $\omega = \frac{\text{Max} - \text{Min}}{K}$)。

特点：适用于数据分布均匀的情况，尤其是在没有先验知识的情况下。对数据质量要求高，对离群值敏感。

等频离散化：

将属性的值域[Min,Max]划分成 K 个区间，每个区间所含的数据个数相等(均为 $\frac{N}{K}$ ，其中 N 为取值总数)。

基于聚类：

对需要离散化的连续性特征，采用聚类算法把样本依据该特征划分成相应的簇或者类。在聚类的结果上，基于特定的策略，决定是否对簇进行进一步的分裂或合并，利用自顶向下的策略针对每一个簇继续运行聚类算法，细分成为更小的子簇，或者利用自底向上的策略对相邻的簇进行合并。在最终确定划分簇后，确定切分点及区间个数。

特点：相似的样本能落到相同的区间段内，可以更好地代表原始数据的信息。

基于信息增益：

该方法源自于决策树模型，在建立决策树时，遍历每一个特征，选择熵最小(信息增益最大)的特征作为正式分裂节点。

- ①对连续型特征进行排序。
- ②把特征的每一个取值作为候选分裂节点(切分点)，计算出相应的熵，选择熵最小的取值作为正式的切分点，将原来的区间一分为二。
- ③递归处理第二步中得到的两个新区间段，直到每个区间段内特征的类别一样为止。
- ④合并相邻的，类的熵值为 0 且特征类别相同的区段，重新计算新区间段类的熵值。
- ⑤重复第四步到满足终止条件(决策树的深度或叶子数)。

基于卡方：

卡方检验属于非参数检验范畴，是一种比较两个总体之间是否存在显著性差异的方法：

$\chi^2 = \sum_i \frac{(A_i - E_i)^2}{E_i}$ ，其中 A_i 为落入区间段的样本个数(观察频数)， E_i 为对应的期望频数，在 n 较大的情况下， χ^2 统计量近似服从自由度为 $k-1$ 的 χ^2 分布。

自下而上的卡方离散化——ChiMerge 方法：

- ①将连续型特征的每一个取值看作是一个单独的区间段，并按照数值大小进行排序。
- ②对每对相邻的区间段计算卡方统计量。卡方值最小或低于设定阈值的相邻区间段合并。
- ③对于新的区间段，递归进行第 1,2 步，直到满足终止条件。

特点：Chimerge 每次迭代的过程只能合并两个区间，如果数据量大，则算法的开销会比较大。需要设定显著性水平，然后计算卡方统计量，通过阈值的设定来控制区间数的多少。

自上而下的卡方离散化：

- ①将连续型特征的所有数据点看作是一个大区间段，并按照数值大小进行排序。
- ②对每对相邻的数据点计算卡方统计量。卡方值最大或高于设定阈值的相邻区间段合并。
- ③对于新的区间段，递归进行第 1,2 步，直到满足终止条件。

卡方值 Chi-square value = $\sum [(观察频数 - 期望频数)^2 / 期望频数]$ ：

$\chi^2 = \sum_{i=1}^2 \sum_{j=1}^C \frac{(A_{ij} - E_{ij})^2}{E_{ij}}$ ，其中 A_{ij} 为第 i 区间段内类别为 j 的样本个数， k 为比较的区间个数，

C 为类别个数， $E_{ij} = \sum_{j=1}^C A_{ij} \cdot \frac{\sum_{i=1}^k A_{ij}}{n}$ 。

类别属性依赖最大化：

类别属性依赖最大化 CAIM(Class-attribute interdependence Maximization)是一种基于熵的特征离散化方法。

与 ChiMerge 不同，它采取自顶向下的策略。通过选择切分点 p ，把特征的取值空间划分为 $f \leq p$ 和 $f > p$ 两个子区间段，用来衡量切分点选择优劣的度量方法是类别属性的相互依赖程度。

假设某个连续型特征有 n 个取值, C 个类别. 假设我们把特征划分为 k 个子区间段, 子区间段集合记为:

$D = \{(d_0, d_1], (d_1, d_2], \dots, (d_{k-1}, d_k]\}$, 其中 d_0 和 d_k 分别为特征的最小值和最大值, $n_{i\cdot}$ 表示属于类别 i 的样本个数, $n_{\cdot j}$ 表示落在区间段 $(d_{j-1}, d_j]$ 的样本个数, n_{ij} 表示在区间 $(d_{j-1}, d_j]$ 内的且属于类别 i 的样本个数。我们可以得到一个由类别特征和离散化特征取值所构成的二维表:

类别	$(d_0, d_1]$	$(d_1, d_2]$	\dots	$(d_{k-1}, d_k]$	类别样本数
1	n_{11}	n_{12}	\dots	n_{1k}	$n_{1\cdot}$
2	n_{21}	n_{22}	\dots	n_{2k}	$n_{2\cdot}$
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
C	n_{C1}	n_{C2}	\dots	n_{Ck}	$n_{C\cdot}$
区间样本数	$n_{\cdot 1}$	$n_{\cdot 2}$	\dots	$n_{\cdot k}$	n

评价当前离散化的好坏的评价准则 $CAIM = \frac{\sum_{j=1}^K \frac{M_j^2}{n_{\cdot j}}}{N}$ 。其中 $M_j = \max\{n_{1j}, n_{2j}, \dots, n_{Cj}\}$ 。

$CAIM$ 的取值为区间 $(0, 1]$ 。 $CAIM$ 的值越大, 说明类和离散区间的相互依赖程度越大, 也就说明了离散化效果越好。

①对进行离散化的特征进行升序排列, 确定取值区间的最小值 d_0 和最大值 d_k , 初始化划分策略 $D = [d_0, d_k]$ 。

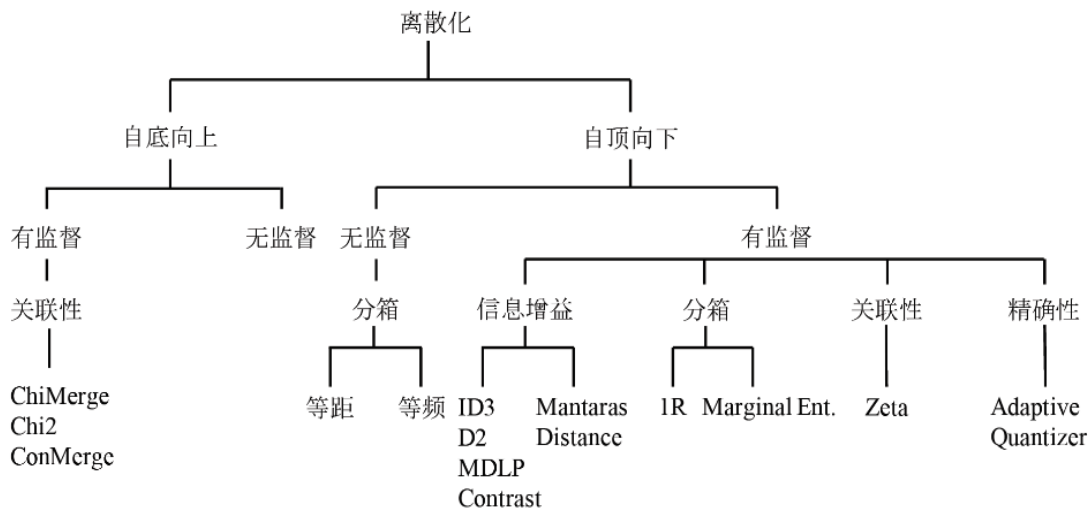
策略 $D = [d_0, d_k]$ 。

②把区间内的每个值当作候选切分点, 计算把区间二分后的 $CAIM$ 值, 并选取 $CAIM$ 值最高的点作为切分点, 并更新 D 。

③对于 D 中的每个区间段, 重复第二步的计算, 直到满足终止条件。

特点: $CAIM$ 只关注区间内拥有最多样本的类别与特征之间的关系, 忽略同一区间内其他类别的信息, $CAIM$ 最终生成的离散化区间个数往往与样本的类别个数接近。

小结:



文献阅读:

【1】Salvador Garcia, Julian Luengo, etc. "A survey of discretization techniques: Taxonomy and empirical analysis in supervised learning". In: IEEE Transaction on knowledge and data engineering. 25(4), 2013, pp. 734-750.

【2】Lukasz A. Kurgan and Krzysztof J. Cios. "CAIM discretization algorithm". In: IEEE transactions on Knowledge and Data Engineering 16.2 (2004), pp. 145-153.

【3】Francis E.H. Tay, "A modified Chi2 Algorithm for Discretization", IEEE Transactions on knowledge and data engineering, 14(3), 2002.

【4】桑雨. 连续数据离散化方法研究.[博士学位论文], 大连理工大学

9. 标称数据的概念分层

便于将标称数据变换到多个粒度层。

以地址这个标称属性集为例，地址这个集合包含{区,市,省,国家}4个标称属性，但没有相关数据语义的知识，未说明其偏序，如何找出任意的标称属性的分层序呢？

启发规则：**较高概念分层的属性，通常包含较少的不同值**。国家这一属性包含的不同值的个数和省或者州包含的不同取值数比较，前者个数显然较少。

回归与分类 模型

先谈机器学习

一、机器学习

1.要素

任务 T(Task): 对数据集 M 分类, 确定给定的某个数据 M_i 的类别。

性能指标 P(Performance): 分类的准确率。

经验来源 E(Experience): 大量的数据以及其对应的类别。

2.应用

数据挖掘: 统计分析、分析评估。

计算机视觉: 图像处理。

自然语言处理 NLP: 文本分类与聚类、信息检索与过滤、机器翻译。

生物特征识别: 人脸识别。

3.术语

数据集 Data Set 里的每条数据称为一个**样本 Sample** 或者**示例 Instance**。样本信息的结果(所属类别)被称为**标记 Label**。

数据对象的**特征 Feature** 和**属性 Attribute** 是一个含义。属性上的值称为**属性值**。

因为一个样本按照其 N 个属性建立 N 维空间, 那么每一个样本都能在该空间中找到自己的坐标向量, 所以样本也被称为**特征向量 Feature Vector**。[注: 此处特征向量不同于特征值分解(eigen decomposition)中的特征向量(eigenvector)]

从数据中学得模型的过程称为**学习**或**训练**, 每个样本被称为**训练样本**, 训练样本组成的数据集称为**训练集**。

分类 Classification: 预测的目标属性是**离散**的。**回归/预测 Regression**: 预测的目标属性是**连续**的。**聚类 Clustering**: 将数据集分成若干个组(簇 Cluster)。

分为**监督学习 Supervised Learning**(分类回归)和**无监督学习**(归纳性学习)Unsupervised Learning(聚类降维)。所谓“监督”, 就是预设的规则/数据的标签。

机器学习 machine learning: 从数据中提取知识, 也被称为预测分析(predictive analytics)或统计学习(statistical learning)。

4.模型评估

模型的目的是为了**泛化 Generalization**, 也就是模型适应新样本的能力。由于使用训练集评估模型会严重高估模型的准确率, 因此需要将数据集 D 划分为训练集 S 和测试集 T 。

(1)留出法 Hold-Out:

正样本 Positive Sample: 在二元分类问题中, 正样本通常代表我们**感兴趣**的类别或事件。例如, 在一个医疗诊断时正样本可以代表患有某种疾病的患者, 而在垃圾邮件过滤器时正样本可以代表垃圾邮件的样本。

负样本 Negative Sample: 在二元分类问题中, 负样本代表与正样本相反类别或事件。例如, 负样本可以代表没有患有疾病的患者, 非垃圾邮件。另外, 比如检测人脸时, 人脸是正样本, 非人脸则是负样本, 比如旁边的树、花之类的其他东西。

在多类别分类问题中, 可以将其中一个类别定义为正样本, 而其他类别被认为是负样本。

留出法将 D 划分为两个互斥的集合 S, T (如 70% S , 30% T), 且对正负样本采用分层抽样。

(2)交叉验证法 Cross Validation(k 折交叉验证法):

将包含 m 个样本的 D 划分为 k 个大小相似的互斥子集，每个子集 $D_i(i=1\sim k)$ 均通过分层抽样得到，使用 $k-1$ 个子集的并集作为 S ，余下的一个子集作为 T 。依次方法进行 k 次训练和测试，最终返回值是这 k 次测试结果的均值。

当 $k=1$ ，交叉验证法退化为留出法。当 $k=m$ ，每个子集只有一个样本数据，就称为留一法 Leave-One-Out(LOO)，评估准确但训练 m 个模型的开销大。因此一般取 $k=5\sim 10$ 。

(3)自助法 Bootstrap:

每次随机从 D 中挑选一个样本复制到 D' 中，重复 m 次。得到的 D' 包含 m 个样本，但约有 $36.8\%(1/e)$ 的样本未在 D' 中出现，这部分作为 T ， D' 中的数据作为 S 。

方法	采样方法	和原始数据分布	相比原始数据容量	适用数据集	估计偏差
留出	分层抽样	不相同	变小	大	存在
交叉	分层抽样	不相同	变小	大	存在
自助	放回抽样	不相同	不变	小	存在

5.性能度量

对于一个机器学习预测模型 f ，预测任务的性能指标可以通过一个函数 $P(D, f)$ 来表示，那么机器学习的过程则是在一个给定的模型空间 \mathcal{F} 中，寻找可以最大化性能指标的预测模型 f^* ，即： $f^* = \arg \max_{f \in \mathcal{F}} P(D, f)$ 。

有监督：降低预测偏差， $f^* = \arg \min_{f \in \mathcal{F}} \frac{1}{|D|} \sum_{(x,y) \in D} L(y, f(x))$

无监督：最大化似然， $p^* = \arg \max_p \frac{1}{|D|} \sum_{x \in D} \log p(x)$

泛化能力：在给定数据分布下的损失函数值的期望， $R(f) = \int_{(x,y)} p(x,y) L(y, f(x)) d(x,y)$

损失函数 Loss Function：单个样本，算的是一个样本的误差。

代价函数 Cost Function：整个训练集，是所有样本误差的平均，也就是损失函数的平均。

目标函数 Object Function：Cost Function + 正则化项。

(1)在分类中：

混淆矩阵 Confusion Matrix:

真实数据 \ 预测结果	正样本	负样本
正样本	TP	FN
负样本	FP	TN

真阳 True Positive: 被预测为正的样本。 假阴 False Negative: 被预测成负的正样本。

假阳 False Positive: 被预测为正的负样本。 真阴 True Negative: 被预测成负的负样本。

精度/准确率 Accuracy: $A = \frac{TP + TN}{TP + FN + FP + TN}$ = 正确分类数/全部分类数。

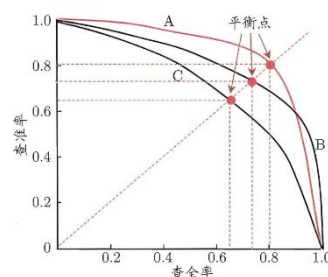
错误率 Error Rate: $E = \frac{FP + FN}{TP + FN + FP + TN}$ = 错误分类数/全部分类数。

查准率/精确率 Precision: $P = \frac{TP}{TP + FP}$ ，分类成正样本的确实是正样本的正确率。

查全率/召回率 Recall: $R = \frac{TP}{TP + FN}$ ，将所有的正样本正确找出来的比例。

显然 P 和 R 不可兼得，如商品推荐算法更多考虑准，逃犯信息检索更多考虑全。P-R 曲线如右图。还有平衡点、ROC 曲线、AUC 等指标。

以信息检索为例，刚开始在页面上显示的信息是用户可能最感兴趣的信息，此时查准率高，但只显示了部分数据，所以查全率低。随着用户不断地下拉，信息符合用户兴趣的匹配程度逐渐降低，查准率不断下降，查全率逐渐上升。当下拉到信息底部时，此时的信息是最不符合用户兴趣，因此查准率最低，但所有的信息都已经展示，查全率最高。



F1 度量: $F1 = \frac{2PR}{P+R}$, 即加权调和平均 $\frac{2}{\frac{1}{P} + \frac{1}{R}}$, 这使得 F1 接近于 P,R 中更小的那个,

所以当 P,R 接近时 F1 值最大。加权调和平均的例子如去程速度 30 返程速度 10 求平均速度。

F-beta 度量: $F_\beta = (1 + \beta^2) \frac{PR}{\beta^2 P + R}$, P,R 权重分别是 $\frac{1}{1 + \beta^2}, \frac{\beta^2}{1 + \beta^2}$, β 越小 P 越重要。

ROC 曲线(受试者操作特征 receiver operating characteristic):

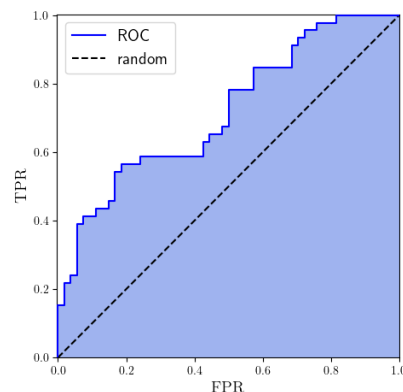
假阳性率 $FPR = \frac{FP}{FP + TN}$, 真阳性率(召回率) $TPR = \frac{TP}{TP + FN}$ 。

ROC 曲线从右上角到左下角，阈值由 0 增大到 1。

例如某样本真实类别为负类，在阈值增大时，其他样本的类别预测没有变化，仅有这一个样本的预测由正类变为负类，那么假阳性减少，真阴性增加，从而 TPR 不变、FPR 减小，是水平线。

如果该样本真实类别是正类，那么真阳性减少、假阳性不变，从而 TPR 减小，FPR 不变，是竖直线。

仅当不少于两个且类别不同的样本被模型赋予了完全相同预测值，ROC 曲线才会出现斜线的。



AUC: 是 ROC 曲线与 x 轴和直线 $x=1$ 围成的面积，称为曲线下面积 area under the curve。

如果将样本按模型预测的正类概率值从小到大排序，AUC 的值等于负类排在正类前面的概率。当正类的预测值都比负类大的时候，AUC=1，也就是完美分割。

基于概率度量的损失函数:

记 p_{y_i} 是 y_i 的概率预测(\hat{y}_i)，有 $p_{y_i} = \frac{\exp(o_i)}{\sum_k \exp(o_k)}$ (o_i 是输出层中的类别 i)。

在二分类中， $y_i = 0, 1$ 分别代表第 i 个样本属于负类/正类。在多分类中， $y_{i,k} = 0, 1$ 分别代表第 i 个样本不属于/属于类别 k 。

KL 散度函数(相对熵)(Kullback-Leibler divergence): $KL = \sum_{i=1}^n y_i \ln \frac{y_i}{p_{y_i}}$

交叉熵 Cross-Entropy: $CE = -\sum_{i=1}^n y_i \ln p_{y_i}$ 。

二分类 $CE = -\frac{1}{n} \sum_{i=1}^n [y_i \ln p_{y_i} + (1 - y_i) \ln (1 - p_{y_i})]$, k 分类: $CE = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K y_{i,k} \ln(p_{y_{i,k}})$ 。可参考 BV15V411W7VB

平衡交叉熵函数 balanced cross entropy: 原本二分类 $CE = \frac{1}{n} \left[\sum_{i=1}^T -\ln p_{y_i} + \sum_{i=1}^F -\ln (1 - p_{y_i}) \right]$,

当正类别 T 远多于负类别 F 时，会使得模型对负类别的性能较差。于是引入权重参数 α ，即

$$\text{BCE} = \frac{1}{n} \left[\sum_{i=1}^T -\alpha \ln p_{y_i} + \sum_{i=1}^F -(1-\alpha) \ln (1-p_{y_i}) \right], \text{ 其中 } \frac{\alpha}{1-\alpha} = \frac{n}{m} \Rightarrow \alpha = \frac{n}{m+n} \in [0,1]。$$

$$\text{focal loss: FL} = \begin{cases} -(1-p_{y_i})^\gamma \ln(p_{y_i}) & \text{if } y=1 \\ -p_{y_i}^\gamma \ln(1-p_{y_i}) & \text{if } y=0 \end{cases} \xrightarrow{p_t = \begin{cases} p_{y_i} & \text{if } y=1 \\ 1-p_{y_i} & \text{if } y=0 \end{cases}} -(1-p_t)^\gamma \ln p_t, \quad p_t \text{ 越大说明越易分, FL 相当于增加了难分样本 } (p_t \text{ 小}) \text{ 在损失函数的权重。}$$

明越易分, FL 相当于增加了难分样本 (p_t 小) 在损失函数的权重。

Hinge Loss 合页损失: $\text{HL}(y_i, f(x_i)) = \max\{0, 1 - y_i \hat{y}_i\}$ 常用于二分类, 此处 \hat{y}_i 取 ± 1 代表正负类。 y_i, \hat{y}_i 同号 $|\hat{y}_i| \geq 1 \Rightarrow \text{HL} = 0, |\hat{y}_i| < 1 \Rightarrow \text{HL} = 1 - y_i \hat{y}_i \in (0, 1); y_i, \hat{y}_i$ 异号 $\text{HL} = 1 - y_i \hat{y}_i \in (1, +\infty)$ 。三个状态对应预测正确, 预测正确但信心不足, 预测错误。

(2) 在回归中:

设数据集 $D = \{x_1, \dots, x_n\}$ 的预测标签 $\hat{Y} = \{\hat{y}_1, \dots, \hat{y}_n\}$ 对应的数据标签是 $Y = \{y_1, \dots, y_n\}$ 。

$$\text{决定系数 } R^2 = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (\bar{y}_i - y_i)^2}, \text{ 取值范围 } [0, 1], \text{ 越大说明拟合效果越好。但随着样本数量增加, } R^2 \text{ 亦会增加。于是引入校正决定系数 Adjusted } R^2 = 1 - \frac{(1-R^2)(n-1)}{n-p-1}, \text{ 其中 } n \text{ 为样本数量, } p \text{ 为特征数量。}$$

量增加, R^2 亦会增加。于是引入校正决定系数 Adjusted $R^2 = 1 - \frac{(1-R^2)(n-1)}{n-p-1}$, 其中 n 为样本数量, p 为特征数量。

基于距离度量的损失函数:

$$\text{平均绝对误差 (L1 loss) Mean Absolute Error (MAE): } \text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$\text{均方误差 (L2 loss) Mean Squared Error (MSE): } \text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{均方根误差 (RMSE) } \text{RMSE} = \sqrt{\text{MSE}}$$

$$\text{Huber 损失: } \text{Huber} = \begin{cases} |y_i - \hat{y}_i| - \frac{\sigma}{2} & |y_i - \hat{y}_i| > \sigma \\ \frac{1}{2\sigma} (y_i - \hat{y}_i)^2 & \text{else} \end{cases}, \text{ 满足 } \begin{cases} \sigma \rightarrow 0: \text{MAE} \\ \sigma \rightarrow \infty: \text{MSE} \end{cases}$$

$$\text{均方对数误差 (MSLE) } \text{MSLE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\ln y_i - \ln \hat{y}_i)^2}$$

$$\text{Smooth L1 损失函数: } \text{Smooth L1} = \begin{cases} \frac{1}{2} (y_i - \hat{y}_i)^2 & |y_i - \hat{y}_i| < 1 \\ |y_i - \hat{y}_i| - \frac{1}{2} & |y_i - \hat{y}_i| \geq 1 \end{cases}$$

其中:

MAE 对离群点有很好的鲁棒性, 但它在残差为零处却不可导, 且更新的梯度始终相同 (即使很小的损失值, 梯度也很大), 不利于模型的收敛。

MSE 无参数、计算成本低、具有明确物理意义, 但在图像和语音处理方面表现较弱。

RMSE 具有凸性和可微性, 且在独立同分布的高斯噪声情况下, 它能提供最大似然估计, 使得它成为回归问题、模式识别、图像处理中最常使用的损失函数。

Huber 损失 是 MAE, MSE 的综合, 导数连续且梯度随误差减小, 可取得更精确的最小值。对异常点具有更好的鲁棒性。但引入了额外的参数, 而且选择合适的参数比较困难。

$$\text{MSEL 用于评估相对误差 } \frac{y - \hat{y}}{y}, \text{ 因为 } |\ln y - \ln \hat{y}_i| \leq \delta \Rightarrow e^{-\delta} \leq \frac{\hat{y}_i}{y} \leq e^{\delta}。$$

Smooth L1 损失主要用在目标检测中防止梯度爆炸。

(3)在聚类中:

外部指标(参考):

将聚类结果与某个“**参考模型**”(如领域专家的划分结果,比如对数据进行标记)进行比较。默认参考模型的性能指标是对样本的最优划分,度量的目的就是使聚类结果与参考模型尽可能相近,也就是聚类结果中被划分在同一簇样本与参考模型样本也被**同样**划分到一个簇的概率越高越好。

给定数据集 $D = \{x_1, \dots, x_n\}$ 经过聚类划分的簇为 $C = \{C_1, \dots, C_k\}$, 参考模型给出的簇划分为 $C^* = \{C_1^*, \dots, C_s^*\}$ 。设 k, k^* 为数据在 C, C^* 中的簇标记向量, 定义:

$$a = |SS|, SS = \{(x_i, x_j) \mid k_i = k_j, k_i^* = k_j^*, i < j\}, \quad b = |SD|, SD = \{(x_i, x_j) \mid k_i = k_j, k_i^* \neq k_j^*, i < j\},$$

$$c = |DS|, DS = \{(x_i, x_j) \mid k_i \neq k_j, k_i^* = k_j^*, i < j\}, \quad d = |DD|, DD = \{(x_i, x_j) \mid k_i \neq k_j, k_i^* \neq k_j^*, i < j\}。$$

其中 S, D 分别表示数据隶属相同/不同簇。集合 SS 表示在 C 中隶属相同簇并且在 C^* 中仍隶属相同簇的样本对。 a, b, c, d 是对应集合中样本对的个数。由于每对样本只能出现在 4 个集合中的一个, 因此 $a + b + c + d = \frac{n(n-1)}{2}$ 。

Jaccard 系数: $JC = \frac{a}{a+b+c}$, 适合二分类问题。

FM 指数: $FMI = \sqrt{\frac{a}{a+b} \cdot \frac{a}{a+c}}$, 结合了 JC 与 RI 的特点。

Rand 指数: $RI = \frac{2(a+d)}{n(n-1)}$, 适合多分类问题。

以上指标取值范围均为 $[0, 1]$, 取值越大代表聚类效果越好。

内部指标(距离):

直接考察聚类结果, 通过计算簇内样本间距与簇间样本距离来评估模型的性能。也就是用簇内样本间距模拟簇内相似度, 用簇间样本距离模拟簇间相似度。定义:

$$avg(C) = \frac{2}{|C|(|C|-1)} \sum_{1 \leq i < j \leq |C|} dist(x_i, x_j), \quad diam(C) = \max_{1 \leq i < j \leq |C|} dist(x_i, x_j),$$

$$d_{\min}(C_i, C_j) = \min_{x_i \in C_i, x_j \in C_j} dist(x_i, x_j), \quad d_{\text{cen}}(C_i, C_j) = dist(u_i, u_j)。$$

其中 $dist$ 用于计算两个样本之间的距离。 $u_i = \frac{\sum_{1 \leq i \leq |C_i|} x_i}{|C_i|}$ 代表簇 C_i 的中心点。 avg 用于计算

对应簇内样本间的平均距离。 $diam$ 用于计算对应簇内样本间的最远距离。 d_{\min} 用于计算两簇间样本的最短距离。 d_{cen} 用于计算两簇间中心点的距离。

DB 指数: $DBI = \frac{1}{k} \sum_{i=1}^k \max_{j \neq i} \left(\frac{avg(C_i) + avg(C_j)}{d_{\text{cen}}(C_i, C_j)} \right)$

Dunn 指数: $DI = \min_{1 \leq i \leq k} \left\{ \min_{j \neq i} \left(\frac{d_{\min}(C_i, C_j)}{\max_{1 \leq x \leq k} diam(C_x)} \right) \right\}$

前者越小越好, 后者越大越好。

分类与回归

0.说在前头:

分类是一种有监督学习问题,主要目标是**根据已知的类别信息(标签)来训练模型,以预测新数据点的类别**。比如信用风险评估:根据用户历史还款信息预测未来是否违约;电子邮件:根据邮件内容将邮件归类(正常邮件/垃圾邮件)。

决策树 Decision Tree, K 近邻算法 K-Nearest Neighbor(KNN), 朴素贝叶斯 Naive Bayes, 支持向量机 Support Vector Machine(SVM)。

回归目标属性 y 是连续的, **预测数据的趋势**, 回归任务同样需要带有标签的训练数据, 但标签是连续数值。比如依据父母身高预测孩子的。

线性回归 Linear Regression, 非线性回归 Nonlinear Regression, Logistic 回归, Lasso 回归, 岭回归 Ridge Regression。

聚类是一种无监督学习任务,其主要目标是**将数据集中的数据点划分为不同的簇,试图发现数据中的内在结构**。通常只需原始数据,不需要标签信息。比如依据兴趣、行为等属性对社交媒体用户分群。(注意这里“分类”的依据是已知的“属性”,而不是已知的“标签”)。

K-means 聚类 K-means Clustering, 层次聚类 Hierarchical Clustering, 密度聚类 Density-Based Spatial Clustering of Applications with Noise (DBSCAN)。

1.分类-决策树

决策树是一种分类方法,从一组无次序无规则的元组中推理出决策树表示形式。采用自顶而下的递归方式,在决策树内部节点进行属性值的比较,进行向下分支,叶子是要学习划分的类。

常用算法: CLS, ID3, C4.5, CART(Classification And Regression Tree), SLIQ(Super-vised Learning In Quest), SPRINT(Scalable Parallelizable Induction of Decision Trees)。

<https://zhuanlan.zhihu.com/p/517743953>

1966 CLS(学习单个概念)→1979 ID3(选择**信息增益**作为特征选择方法)→1986 ID4(在每个可能的决策树节点创建缓冲区,使决策树可以递增式生成,即在不断增加数据的情况下持续构建决策树)→1988 ID5(提高效率)→1993 C4.5(选择**信息增益率**作为特征选择方法,同时处理连续和离散特征,还能处理缺失值)

另一类决策树算法为 CART,与 C4.5 不同的是, CART 的决策树由二元逻辑问题生成,每个树节点只有两个分枝,分别包括学习实例的正例与反例。CART 这种树构建算法既可用于分类也可用于回归。

递推过程:划分数据集时寻找起决定性作用的特征。

决策树的**特点**: 1.推理过程容易理解,决策推理过程可以表示成 If Then 形式。2.推理过程完全依赖于属性变量的取值特点。3.可能忽略目标变量没有贡献的属性变量,这也为判断属性变量的重要性,减少变量的数目提供参考。

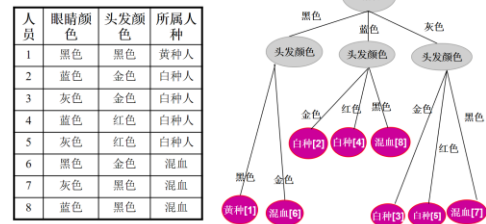
CLS 算法 Concept Learning System

CLS 是早期的决策树学习算法,它是决策树学习算法的基础。

算法思想:

从一棵空决策树开始,选择某一属性(分类属性)作为测试属性,该测试属性对应决策树中的决策结点,根据该属性的值的不同将训练样本分成相应的子集:

如果该子集为空,或该子集中的样本属于同一个类,则该子集为叶结点;否则该子集对应于决策树的内部结点,需要选择一个新的分类属性对该子集进行划分,直到所有的子集都为空或者属于同一类。如右图。



算法问题:

测试属性集的组成以及测试属性的先后对决策树的学习具有举足轻重的影响。

ID3 算法 Iterative Dichotomiser 3

算法思想:

从根节点开始,选择“某”属性特征(如收入)。选择属性特征的策略:**不纯度 impurity**。即落在当前节点的样本类别分布的均衡程度,节点分裂后,节点不纯度应该更低。选择特征及对应分割点,使得分裂前后的不纯度下降最大。

节点不纯度的度量

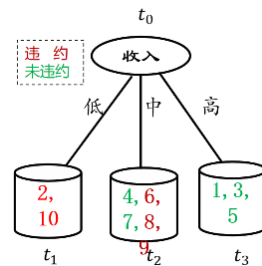
(1)信息熵 Entropy: 描述信息不确定度: $H(x) = -\sum p_i \log_2(p_i)$, 其中 p_i 为 $P(X = x_i) = \frac{N_i}{N}$ 。

当样本均匀分布在每一个类中时,熵为 $\log_2 C$, 说明不纯度大;当所有的样本属于同一个类时,熵为 0, 说明不纯度小。

$$H(t_1) = -\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} = 0$$

例子如右图, $H(t_2) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.971$

$$H(t_3) = -\frac{0}{3} \log_2 \frac{0}{3} - \frac{3}{3} \log_2 \frac{3}{3} = 0$$



条件熵 Condition Entropy: $H(Y|X) = \sum p_i H(Y|X = x_i)$

(2)信息增益 Information Gain: 节点分裂前后信息熵的下降值(得知某特征 A 而使得分类 X 的信息的不确定性减少的程度): $\text{Gain}(A) = H(X) - H(X|A)$ 。其中 $H(X) = H(t_0)$, $H(X|A) =$

$\sum_{k=1}^n \frac{n_k}{n} H(t_k)$, n_k 是子节点 t_k 的样本数量, n 是父节点 t_0 的样本数量。

信息增益越大,说明该特征对结果的影响越大。

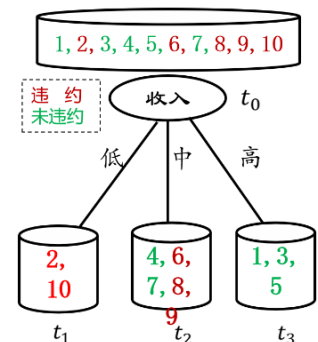
例子如右图, $H(t_0) = -\frac{5}{10} \log_2 \frac{5}{10} - \frac{5}{10} \log_2 \frac{5}{10} = 1$, $H(t_{1,2,3})$ 如前计算,

$$\text{Gain}(\text{收入}) = H(t_0) - [\frac{2}{10} H(t_1) + \frac{5}{10} H(t_2) + \frac{3}{10} H(t_3)] = 0.5145$$

(3)信息增益率 Gain Ratio: $\text{Gain}_{\text{ratio}}(A) = \frac{\text{Gain}(A)}{H(A)}$

(4)基尼指数 Gini Index: $\text{Gini}_{\text{index}}(A) = \sum_{v=1}^n \frac{|X_v|}{|X|} \text{Gini}(X_v)$, 式中 $\text{Gini}(X) = \sum_{k=1}^m p_k(1-p_k) = 1 - \sum_{k=1}^m p_k^2$ 。

$\text{Gini}(X)$ 反应了数据集 X 中随机抽取两个样本,其类别不一致的概率。当样本均匀分布在每



一个类中, $Gini(X) = 1 - \frac{1}{m}$ (m 为该类的样本个数), 说明不纯度大。当样本都分布在同一个类中, $Gini(X) = 0$, 说明不纯度小。

$Gini_{index}(A)$ 就是对于特征 A 下集合的基尼指数, 其中 A 的取值共有 n 种可能, X_v 表示对数据集 X 按照特征 A 划分的数据子集。比如这里特征就是 t_0 (收入), 划分为 $t_{1,2,3}$ 共 $n=3$ 个子集。

$$Gini(t_1) = 1 - [(\frac{2}{2})^2 + (\frac{0}{2})^2] = 0$$

$$Error(t_1) = 1 - \max(\frac{2}{2}, \frac{0}{2}) = 0$$

例子仍同前, $Gini(t_2) = 1 - [(\frac{3}{5})^2 + (\frac{2}{5})^2] = 0.480$

$$Error(t_2) = 1 - \max(\frac{3}{5}, \frac{2}{5}) = 0.4$$

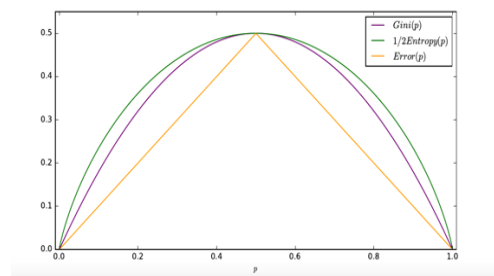
$$Gini(t_3) = 1 - [(\frac{0}{3})^2 + (\frac{3}{3})^2] = 0$$

$$Error(t_3) = 1 - \max(\frac{0}{3}, \frac{3}{3}) = 0$$

ID3 采用信息增益作为特征选择的方法(倾向于分裂成很多小节点「节点的样本数较小」, 容易造成过拟合)。C4.5 采用信息增益率(减少信息增益准则对取值数目较多的属性有所偏好, 避免分裂成过多小节点)。CART 采用基尼指数(选择基尼指数最小的属性进行划分)

(5) 误分率 Misclassification error: $Error(t) = 1 - \max\{p_k\} (k=1, 2, \dots, m)$ 。当样本均匀分布在每一个类中, 误分率为 $1 - \frac{1}{m}$, 不纯度最大。当样本分布在同一个类中, 误分率为 0, 不纯度最小。例子见上文 Gini 后面的 Error。

在同一个类中, 误分率为 0, 不纯度最小。例子见上文 Gini 后面的 Error。



(6) 决策树生成算法的小结:

对于二分类, 如果正类样本的比例为 p , 则有

$$H(t) = -p \log_2 p - (1-p) \log_2 (1-p)$$

$$Gini(t) = 2p(1-p), \text{ 如右图}$$

$$Error(t) = 1 - \max\{p, 1-p\}$$

信息熵最大, 说明对不纯度的惩罚最强。

算法	特征类型	不纯度度量	分割的子节点数量	目标特征类型
ID3	离散型	信息熵(信息增益)	$K \geq 2$	离散型
C4.5	离散型、连续型	信息熵(信息增益率)	$K \geq 2$	离散型
C5.0	离散型、连续型	信息熵(信息增益率)	$K \geq 2$	离散型
CART	离散型、连续型	Gini 指数	$K = 2$	离散型、连续型

决策树生成后转换成 IF-THEN 规则的集合, 使得决策树模型的可解释性较好。

(7) 决策树生成的方法小结:

```
def CreateTree():
    if 该子集为空或该子集中的样本属于同一个类: # 则该子集为叶结点
        return 类标签
    else: # 该子集是内部结点, 需要进行划分。直到所有的子集都为空或属于同一类
        计算信息增益, 确定最好特征。划分数数据集, 创建分支节点。
        for 每个节点:
            CreateTree() 并将返回结果添加到分支节点。
        return 分支节点
```

(8) 决策树的剪枝:

预剪枝：在 CreateTree 时就停止信息量较小[即：信息增益(率)较小]的分支。比如限制深度、叶子节点个数、叶子节点样本数、信息增益量。

后剪枝：先 CreateTree，搞完了再剪枝。通过衡量标准 $C_\alpha(T) = C(T) + \alpha \cdot |T_{\text{leaf}}|$ ，

式中 $C(T)$ 为当前节点的基尼系数， α 自行指定(α 越大越不容易过拟合)，

T_{leaf} 是叶子节点个数。比如节点 $X[3] \leq 1.55$ ，就是比较 $0.4444 \times 6 + 1 \cdot \alpha$ 和

$0.0 \times 3 + 0.4444 \times 3 + 2\alpha$ 的大小。

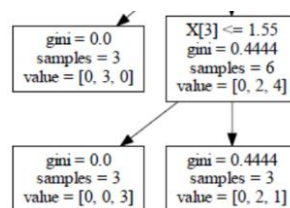
(9)决策树的优缺点：

优点：直观。允许存在缺失值。对离群点不敏感。可用于复杂的非线性关系。

缺点：决策树由于使用贪婪算法的思想，即在每次分裂时选择当前情况下带来信息增益最高的属性进行分裂，容易陷入**局部最优解**。如果缺少合适的剪枝策略或终止分支处理，容易过拟合，比如以 ID 来决策，但无效。

(10)补充：德国信用数据集：

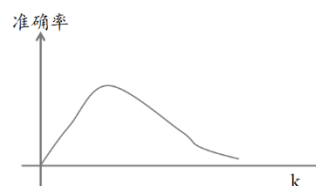
<https://archive.ics.uci.edu/ml/datasets/South+German+Credit+%28UPDATE%29>



2.分类-k 邻近算法(KNN)

当对测试样本进行分类时，找到训练集中与该样本集最相似的 k 个样本，根据 k 个样本的标签(选取其中最多的)确定测试样本的标签。

准确率与 k 的函数一般成单峰。二分类问题中 k 一般设为



奇数，防止出现两种类型出现次数相同的情况。 k 一般小于 30。

寻找相似样本可以认为距离近的数据相似。

算法流程：

- 1.确定 k 的大小和距离计算方法。
- 2.从训练样本中得到 k 个与测试样本最近的样本。
- 3.根据 k 个最相似的训练样本的类别，通过投票的方式来确定测试样本的类别。

距离度量：

定义：

- 1.非负性： $f(A, B) \geq 0$ ，当且仅当 $A = B$ 时等号成立。
- 2.反身性： $f(A, B) = f(B, A)$
- 3.三角不等式： $f(A, B) \leq f(A, C) + f(B, C)$

通用公式：

首先定义 x, y 为任意样本数据点,查询点，有时候也将查询点记为 p 。

明可夫斯基 Minkowski 距离(闵氏距离)： $d(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}$

曼哈顿距离： $p = 1$ 。欧氏距离： $p = 2$ 。切比雪夫距离： $p = +\infty$

曼哈顿距离(城市街区距离): $d(x, y) = \|x - y\|_1 = \sum_{i=1}^D |x_i - y_i|$, 因为开车不能穿墙, 只能沿“水平竖直”道路。

欧氏距离: n 维空间有: $d(x, y) = \|x - y\| = \sqrt{(x - y)^T (x - y)} = \sqrt{\sum_{i=1}^D (x_i - y_i)^2}$

标准化欧氏距离: $d(x, y) = \sqrt{(x - y)^T D^{-1} (x - y)}$, 其中 D 是对角方阵, 对角线上元素为每个特征的方差, 即 $D = \text{diag}(\text{diag}(\Sigma))^{\frac{1}{2}}$ 。其本质是 Z 分数 $z_i = \frac{x_i - y_i}{\sigma_i}$, 证明如下:

$$D = \text{diag}(\text{diag}(\Sigma))^{\frac{1}{2}} = \text{diag} \left(\text{diag} \begin{bmatrix} \sigma_1^2 & \rho_{1,2}\sigma_1\sigma_2 & \cdots & \rho_{1,D}\sigma_1\sigma_D \\ \rho_{1,2}\sigma_1\sigma_2 & \sigma_2^2 & \cdots & \rho_{2,D}\sigma_2\sigma_D \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{1,D}\sigma_1\sigma_D & \rho_{2,D}\sigma_2\sigma_D & \cdots & \sigma_D^2 \end{bmatrix} \right)^{\frac{1}{2}} = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_D \end{bmatrix}$$

$$d(x, y) = \sqrt{(x - y)^T \begin{bmatrix} \sigma_1^2 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \sigma_D^2 \end{bmatrix}^{-1} (x - y)} = \sqrt{\frac{(x_1 - y_1)^2}{\sigma_1^2} + \cdots + \frac{(x_D - y_D)^2}{\sigma_D^2}} = \sqrt{\sum_{i=1}^D \frac{(x_i - y_i)^2}{\sigma_i^2}} = \sqrt{\sum_{i=1}^D z_i^2}$$

马氏 Mahalanobis 距离: $d(x, y) = \sqrt{(x - y)^T \Sigma^{-1} (x - y)}$, 其中 Σ 为协方差矩阵。

如何理解 Σ^{-1} ?

1. 正交变换

协方差矩阵 Σ 半正定, 如果其正定, 则其可对角化 $V^{-1} \Sigma V = \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, 其中 V 是正交矩阵。令 $y = V^{-1}x$, 则 $\text{cov}(y, y) = \text{cov}(V^{-1}x, V^{-1}x) \stackrel{(\text{by Lemma 1})}{=} V^{-1} \text{cov}(x, x) V = \Lambda$ 。注意到此时协方差矩阵 $\text{cov}(y, y)$ 是对角阵, 说明只有每个维度和自己的方差存在, 而不同维度之间的协方差为 0, 故相关系数也为 0, 因此不同维度不相关。每个特征值 $\lambda_i = \text{cov}(y_i, y_i)$, 即方差, 即标准差的平方。

正交变换的作用是: 消除新的随机向量 y 不同维度之间的相关性^[注 1], 同时保持各个样本点的欧式距离保持不变^[注 2]。

引理 1: 首先, 注意到 V^{-1} 是常数矩阵, 有 $E[V^{-1}x] = V^{-1}E[x]$, 因此:

$$\begin{aligned} \text{cov}(V^{-1}x, V^{-1}x) &= E[(V^{-1}x - V^{-1}\mu_x)(V^{-1}x - V^{-1}\mu_x)^T] = E[V^{-1}(x - \mu_x)(x - \mu_x)^T(V^{-1})^T] \\ &= V^{-1}E[(x - \mu_x)(x - \mu_x)^T](V^{-1})^T = V^{-1} \text{cov}(x, x)(V^{-1})^T = V^{-1} \text{cov}(x, x) V \end{aligned}$$

注 1: 做正交变换的目的并不是使不同维度线性无关, 而是使不同维度不相关, 即不同维度的协方差为 0, 相关系数为 0。线性无关=不相关+非线性相关。

注 2: 正交变换无非是旋转或镜像(视行列式的值而定), 自然保持距离不变。

2. 标准化

考虑 $z_i = \frac{x_i - y_i}{\sigma_i}$, 将 σ_i 替换为 1 中的特征值, 即 $z_i = \frac{x_i - y_i}{\sqrt{\lambda_i}}$ 。于是变成 $z = \Lambda^{-\frac{1}{2}} V^{-1} x$, 其

中 $\Lambda^{-\frac{1}{2}} = \text{diag}(\frac{1}{\sqrt{\lambda_1}}, \frac{1}{\sqrt{\lambda_2}}, \dots, \frac{1}{\sqrt{\lambda_n}})$ 。显然 $\text{cov}(z, z) = \Lambda^{-\frac{1}{2}} \text{cov}(y, y) \Lambda^{-\frac{1}{2}} = \Lambda^{-\frac{1}{2}} \Lambda \Lambda^{-\frac{1}{2}} = I$ 。因此对于标准化后的随机向量, 不同维度的协方差仍然都为 0, 而同维度的方差都为 1。

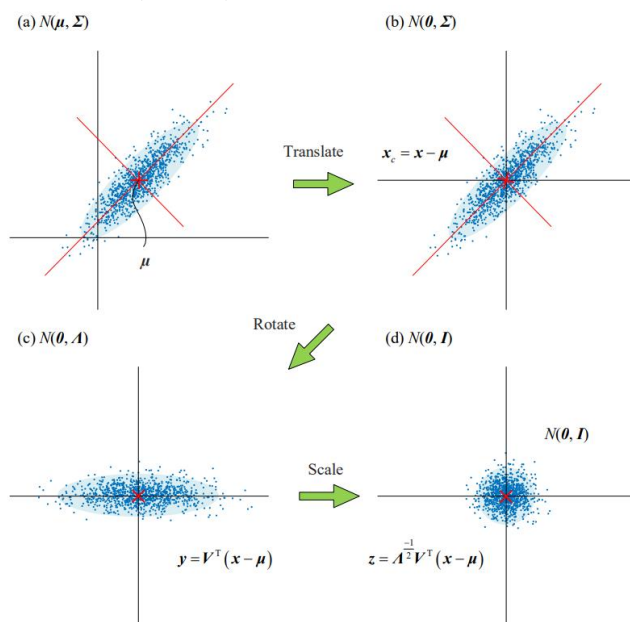
3. 小结(以正态分布中的马氏距离来总结)

根据上文, 易得马氏距离的单位是标准差。比如, 马氏距离计算结果为 3, 应该称作 3 个标准差。因为马氏距离是对原始数据的中心化、旋转、缩放, 反过来, 用马氏距离度量的样本就会自动呈现出以均值为原点, 标准差为单位长度的形式, 公式与图形如下所示:

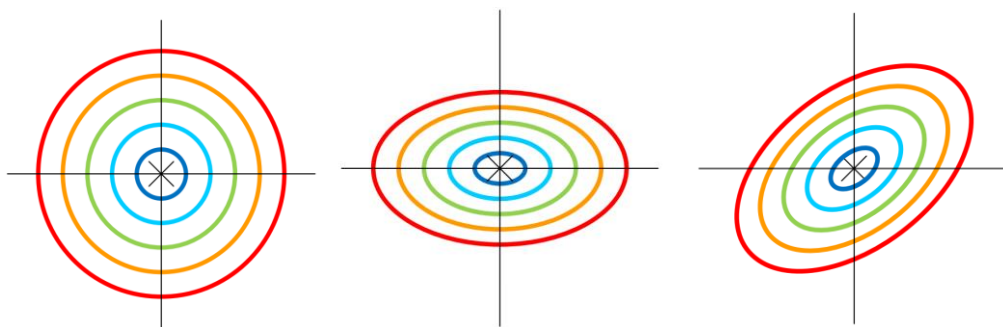
对于 $f_{\mathcal{X}}(\mathbf{x}) = \frac{\exp\left[-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right]}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}}$, 同理 $\boldsymbol{\Sigma} = \mathbf{V} \boldsymbol{\Lambda} \mathbf{V}^T \Rightarrow \boldsymbol{\Sigma}^{-1} = \mathbf{V} \boldsymbol{\Lambda}^{-1} \mathbf{V}^T$ 。那么:

$$(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}) = (\mathbf{x}-\boldsymbol{\mu})^T \mathbf{V} \boldsymbol{\Lambda}^{-1} \mathbf{V}^T (\mathbf{x}-\boldsymbol{\mu}) = \left[\boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{V}^T (\mathbf{x}-\boldsymbol{\mu}) \right]^T \left[\boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{V}^T (\mathbf{x}-\boldsymbol{\mu}) \right] = \left\| \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{V}^T (\mathbf{x}-\boldsymbol{\mu}) \right\|_2^2$$

说明这个过程是将 \mathbf{x} (超椭圆分布) 通过平移 $\mathbf{x}-\boldsymbol{\mu}$ 、旋转/镜像 \mathbf{V}^T 、缩放 $\boldsymbol{\Lambda}^{-\frac{1}{2}}$ 变成了 \mathbf{x}' (超球分布), 如下图是二维上的椭圆变为圆的过程。



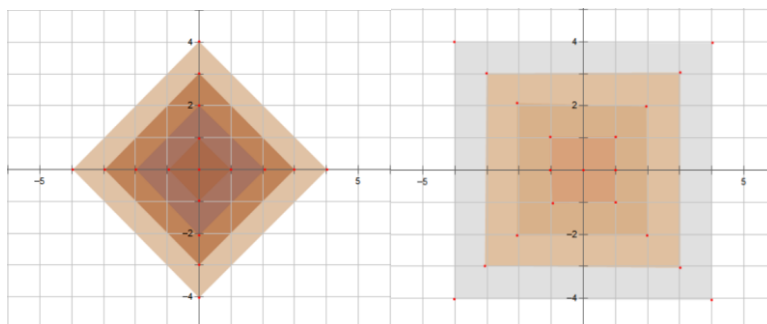
欧式、标准化欧式、马氏距离图像如下:



切比雪夫距离: $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_{\infty} = \max\{|x_i - y_i|\}$

二维时, 曼哈顿下的点 (x, y) 对应于切比雪夫下的 $(x+y, x-y)$ 。

曼哈顿距离和切比雪夫距离都是正方形, 区别在于, 曼哈顿距离是上下左右 4 个方向的延伸, 切比雪夫考虑是四周 8 个方向的延伸。因此曼哈顿是 45° 斜着的正方形, 切比雪夫是正着的正方形。如下图:



亲近度 affinity/相似度 similarity:

亲近度是与距离相反的度量。两个样本数据距离越远，两者亲近度越低。

余弦距离 Cosine Distance: 余弦相似度是类比 $\cos \theta$ 的计算公式定义的。相比距离度量，余弦相似度更加注重两个向量在**方向**上的差异，而非距离或长度上。

$$\text{余弦相似度 } k(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}} \in [-1, 1]$$

余弦距离 $d(\mathbf{x}, \mathbf{y}) = 1 - k(\mathbf{x}, \mathbf{y}) \in [0, 2]$ 值越小越正相似。

相关系数相似度 correlation similarity: $k(\mathbf{x}, \mathbf{y}) = \frac{(\mathbf{x} - \bar{\mathbf{x}})^T (\mathbf{y} - \bar{\mathbf{y}})}{\|\mathbf{x} - \bar{\mathbf{x}}\| \|\mathbf{y} - \bar{\mathbf{y}}\|}$ 相比于余弦相似度，增加了中

心化。

杰卡德距离:

$d_J(A, B)$ 用来衡量两个集合差异性，是杰卡德相似系数的补集，定义为 $1 - J(A, B)$ 。而杰

卡德相似系数(Jaccard similarity coefficient)，也称杰卡德指数(Jaccard Index)，是用来衡量两个集合相似度的一种指标。

$$J(A, B) = \frac{|\#(A \cap B)|}{|\#(A \cup B)|}, \quad d_J(A, B) = 1 - J(A, B) = \frac{|\#(A \cup B)| - |\#(A \cap B)|}{|\#(A \cup B)|}$$

例如：两个 n 维向量样本 A(0111)和 B(1011)。有 $J(A, B) = \frac{p}{p+q+r} = \frac{2}{2+1+1} = 0.5$

p : 样本 A 与 B 都是 1 的维度的个数 q : A1, B 0 r : A0, B1 s : A0, B0。

$$\text{广义 Jaccard 相似系数: } J_g(a, b) = \frac{\sum_i \min(a_i, b_i)}{\sum_i \max(a_i, b_i)}$$

核函数亲近度:

线性核/内积核: $\kappa(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y} = \mathbf{x} \cdot \mathbf{y}$

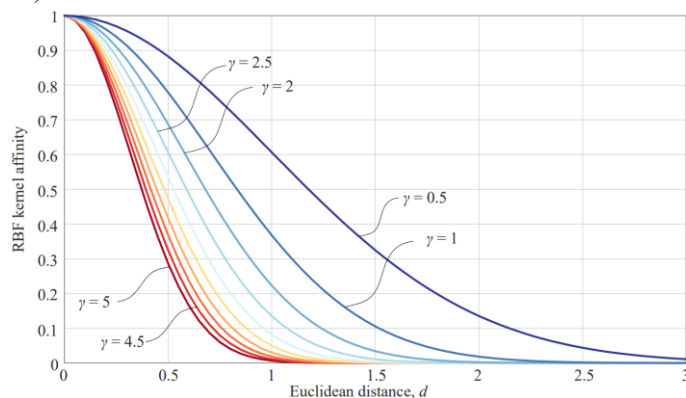
多项式核: $\kappa(\mathbf{x}, \mathbf{y}) = (\gamma \mathbf{x}^T \mathbf{y} + r)^d = (\gamma \mathbf{x} \cdot \mathbf{y} + r)^d$

余弦相似度核: $\kappa(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \in [-1, 1]$

Sigmoid 核: $\kappa(\mathbf{x}, \mathbf{y}) = \tanh(\gamma \mathbf{x}^T \mathbf{y} + r) = \tanh(\gamma \mathbf{x} \cdot \mathbf{y} + r)$ ，带有 Sigmoid 核的 SVM 等价于两层的多层感知机。

拉普拉斯核: $\kappa(\mathbf{x}, \mathbf{y}) = \exp(-\gamma \|\mathbf{x} - \mathbf{y}\|_1)$

高斯核: $\kappa(x, y) = \exp(-\gamma \|x - y\|^2) = \exp(-\gamma d^2) \in (0, 1]$, 距离越小亲近度越高。从距离→亲近度的转换角度来看, 多元高斯分布的分子 $\exp\left[-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right]$ 完成的的就是马氏距离 d 到概率密度(亲近度)的转化。



k 邻近算法小结:

对异常数据不敏感, 具有较好的抗噪性。K 近邻算法的计算效率不高。当训练集较小的时候, K 近邻算法易导致过度拟合。

3.分类-朴素贝叶斯

“朴素”: 假设变量间相互独立(没有这个前提, 朴素贝叶斯模型的分类效果会较差)

$$\Rightarrow P(AB) = P(A)P(B), \text{ 可得: } P(AB|X) = \frac{P(ABX)}{P(X)} = \frac{P(AX)P(BX)}{P(X)P(X)} = P(A|X)P(B|X)$$

贝叶斯公式: $P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$ 。如果用 x 代表属性集合(数据集), y 表示类别标签

(预测), 即 $P(y_i|x) = \frac{P(x|y_i)P(y_i)}{P(x)}$ 。当各个属性相互独立时, 有 $P(x|y_i)P(y_i) = P(y_i) \prod_{j=1}^m P(a_j|y_i)$,

其中 $x = \{a_1, a_2, \dots, a_n\}$, 每个 a_i 代表数据的一个属性, 可以理解为 $x = a_i$ 。

以 $y = \text{True} / \text{False}$ 为例子, 有:

$$P(T|X) = \frac{P(x_1|T)P(x_2|T) \cdots P(x_n|T)P(T)}{P(x_1)P(x_2) \cdots P(x_n)}, P(F|X) = \frac{P(x_1|F)P(x_2|F) \cdots P(x_n|F)P(F)}{P(x_1)P(x_2) \cdots P(x_n)},$$

不难发现分母相同, 于是比较分子大小即可完成分类(分类为 T 或 F)。推广到一般, 就是当分子 $P(y_i) \prod_{j=1}^m P(a_j|y_i)$ 最大时, 预测 x_i 最可能标签为 y_i 。

该思想的来源:

比如要求 $P(\text{违约}|\text{男}, \text{收入中}, \text{本科}, \text{未婚})$ 的值大多数情况下很难, 但 $P(\text{男}|\text{违约}), P(\text{收入中}|\text{违约}), P(\text{本科}|\text{违约}), P(\text{未婚}|\text{违约})$ 一般是容易的。则有:

$$P(\text{违约}|\text{男}, \text{收入中}, \text{本科}, \text{未婚}) = [P(\text{男}|\text{违约}) \cdot P(\text{收入中}|\text{违约}) \cdot P(\text{本科}|\text{违约}) \cdot P(\text{未婚}|\text{违约}) \cdot P(\text{违约})] / [P(\text{男}) \cdot P(\text{收入中}) \cdot P(\text{本科}) \cdot P(\text{未婚})]$$

0 概率处理: 可能某个特征属性在训练集中没有出现, 这样如果按照以上的计算方法, 就会导致后验概率为 0 的情况, 所以要对条件概率的求解公式修正: $P(B_i) = \frac{0+1}{N}$ (加 1 平滑)

拉普拉斯平滑: $P(X = x_i | Y = k) = \prod_{i=1}^d \frac{\sum_{j=1}^n I(X_i^{(j)} = x_{ij}, y_i = k) + \alpha}{\sum_{j=1}^n I(y_i = k) + \lambda \alpha}$ 。 λ 为特征的维数, α 是平滑值。

x_{ij} 代表第 i 个特征的第 j 个选择。 I 指示函数(indicator function)用于描述某个事件是否发生的函数

$$I(y_i = k) = \begin{cases} 1, & \text{if } y_i = k \\ 0, & \text{if } y_i \neq k \end{cases}$$

比如特征为有钱/没钱, 维数 $\lambda = 2$ 。 $\alpha = 1$ 时称为 Laplace 平滑, $0 < \alpha < 1$ 时称为 Lidstone 平滑, $\alpha = 0$ 说明不做平滑。

当特征是连续变量时,

1. 将连续的特征离散化。困难: 如何离散化? (tip: 等频分箱[数据量相同]、等宽分箱[宽度相同]、基于聚类的分箱、基于决策树的分箱)

2. 采用高斯模型, 假定连续型特征服从正态分布, 使用正态分布对条件参数进行求解。

$$P(X = x_i | y = c) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

朴素贝叶斯算法**优缺点分析:**

分类效率良好, 逻辑简单, 易于实现。对小规模数据表现良好, 能处理多分类任务, 适合增量式训练。特征之间相互独立, 只会涉及二维存储, 分类算法过程开销小。对缺失数据不敏感, 常用于文本分类。

在属性个数多且属性之间的相关性比较大时, 分类效果不理想。分类过程中需要知道先验概率和条件参数, 而先验概率和条件参数的求解过程取决于我们选择的模型, 当选择不合适的模型时, 求得的先验概率和条件参数的值也会影响分类的结果。

例子: 将邮件划分为正常邮件和垃圾邮件:

不妨记正常邮件为 T , 垃圾邮件为 F 。训练集邮件总数为 N , 训练集正常邮件总数为 N_T , 训练集垃圾邮件总数为 N_F , 某单词 C 出现在正常或垃圾邮件的数目分别为 C_T 与 C_F , 有:

$$P(T) = \frac{N_T + 1}{N + 2}, P(F) = \frac{N_F + 1}{N + 2}, P(C|T) = \frac{C_T + 1}{N_T + 2}, P(C|F) = \frac{C_F + 1}{N_F + 2}.$$

4. 分类- SVM 支持向量机 Support Vector Machine

解决小样本、非线性及高维模式识别。根据有限的样本信息在模型的复杂性(即对特定训练样本的学习精度)和学习能力(即无错误地识别任意样本的能力)之间寻求最佳折衷(全局最优解), 以期获得最好的泛化能力(通过使用最大分类间隙来设计决策最优分类超平面)。

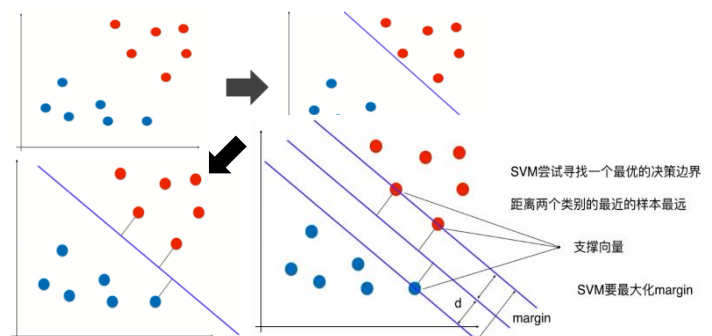
应用: 在模式识别领域中的文本识别与分类、人脸识别等, 工程技术、信息过滤。

寻找决策边界:

决策边界明显不唯一, 需要找更“好”的。量化标准就是“分类间隔”(两虚线的间隔)。

分割线的特点: 找到红蓝点中离这根直线最近的点, 同时尽可能使得这些点离这根直线尽可能的远。总的来说就是这根直线离蓝色的样本尽可能远, 离红色样本也尽可能远。即: 使两虚线中的最小距离取得最大值(不过图四中的两虚线画成实线了, 「支持向量所在的两根线」)。

那如何求这个距离 d : 高维(n 维欧式空间)中某个点 $M_0: x_0 = (x_0^{(1)}, x_0^{(2)}, \dots, x_0^{(n)})$ 到某一个超



平面 $L: \omega \cdot x + b = 0$ (法向量 $\omega = (\omega^{(1)}, \omega^{(2)}, \dots, \omega^{(n)})$) 的距离 $d = \frac{|\omega \cdot x_0 + b|}{\|\omega\|}$ 。

想要距离最大，可证就要法向量长度 $\|\omega\|$ 取最小值，也就是 $\min \frac{1}{2} \|\omega\|^2$ ，证明如下：

样本特征向量记为 x_i ，类别标签 $y_i = 1, -1$ 代表正/负样本。对正/负样本分别有 $\begin{cases} \omega^T x + b > 0 \\ \omega^T x + b < 0 \end{cases}$ ，

可以统一写为 $y_i \cdot (\omega^T x_i + b) > 0$ 。因为超平面方程两边能同乘不为 0 的常数，不妨为 ω, b 加上约束： $\min_x |\omega^T x_i + b| = 1$ ，于是判断超平面是否正确分类的依据变为 $y_i \cdot (\omega^T x_i + b) \geq 1$ 。于是分

类间隔是： $d(\omega, b) = \min_{x_i, y_i=1} \frac{|\omega^T x_i + b|}{\|\omega\|} + \min_{x_i, y_i=-1} \frac{|\omega^T x_i + b|}{\|\omega\|} = \frac{\min_{x_i, y_i=1} |\omega^T x_i + b| + \min_{x_i, y_i=-1} |\omega^T x_i + b|}{\|\omega\|} = \frac{2}{\|\omega\|}$ 。要使

间隔最大化，就是 $\min \frac{1}{2} \|\omega\|^2$ 。于是问题转化为：

$$\begin{aligned} & \min \frac{1}{2} \omega^T \omega \\ & s.t. \quad y_i \cdot (\omega^T x_i + b) \geq 1 \end{aligned}$$

为提高容错能力(泛化能力)，允许一些点出现在虚线和支撑向量所在直线(两条直线)之间。因此我们在不等式右边减去一个整数 ζ (zeta)(罚分值 Penalty)，也就是 $y_i \cdot (\omega^T x_i + b) \geq 1 - \zeta_i$ ，

于是要求的是 $\min \frac{1}{2} \|\omega\|^2 + \sum_{i=1}^m \xi_i$ ，但是显然 ζ 不能很大，为此我们还要限定 ζ 的大小：

$\min \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m \xi_i$ (超参数 C 越大容错空间越小， C 越小容错空间越大)。此时称为 Soft Margin SVM(之前的称为 Hard Margin SVM)。最终就要不断最小化罚分总和。

可以化为：约束条件： $\sum_{i=1}^N \alpha_i y_i = 0, 0 \leq \alpha_i \leq C (i=1, 2, \dots, N)$

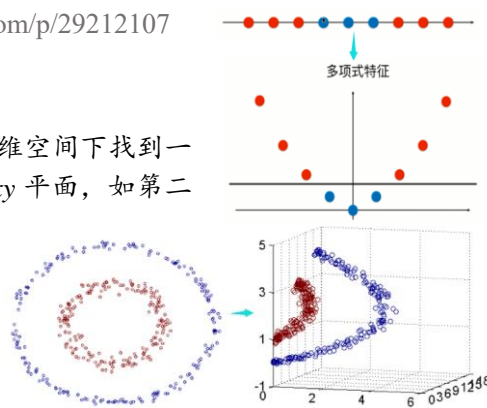
$$\text{目标函数: } \min \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i^T x_j) - \sum_{i=1}^N \alpha_i$$

具体可参考序列最小优化 [SMO 算法](https://zhuanlan.zhihu.com/p/29212107) <https://zhuanlan.zhihu.com/p/29212107>

线性不可分：升维

因为升高维度，数据会显得“稀疏”。如右图，没办法在一维空间下找到一个点分割红蓝样本。现在选择升高维度，使得样本能够坐落在 xy 平面，如第二个特征值设置为 x^2 。这些点在 x 方向上的位置并没有改变，但是因为加入了 y 这个维度，此时就能找到一根直线分割。

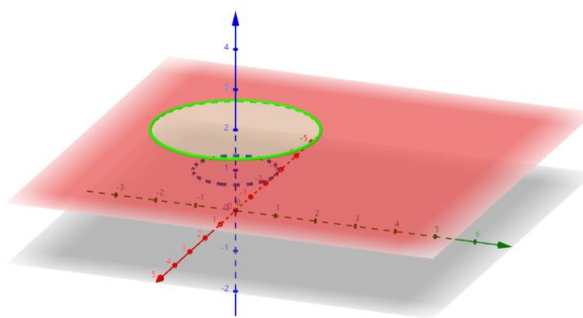
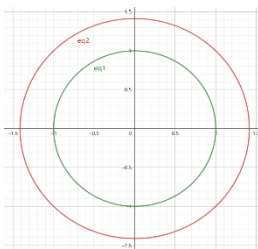
核函数 Kernel 的作用：从低维空间到高维空间的映射，而这个映射可以把低维空间中线性不可分的两类点变成线性可分的(如右图)，其目的与多项式特征中升维是一样的道理。



升维的一个例子：平面上无法用直线区分

$$x^2 + y^2 = 1 \text{ 与 } x^2 + y^2 = 2,$$

$$\text{但 } (x, y) \rightarrow (x, y, \sqrt{x^2 + y^2})$$

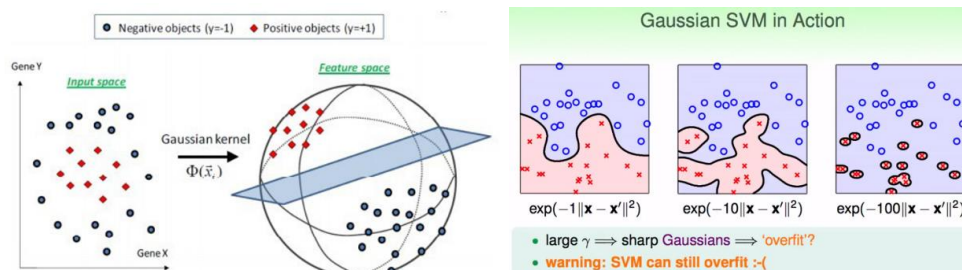


后在三维空间中可由 $z = \frac{3}{2}$ 区分这两个平面。

高斯核函数 RBF(径向基核函数): $K(x, y) = \exp(-\gamma \|x - y\|^2) \in (0, 1]$, 当 $x = y$ 时函数值为 1。

γ 越大, 二元高斯核函数 $\exp(-\gamma(x_1^2 + x_2^2))$ 坡面越陡峭。RBF 也可写成 $K(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$ 。

RBF 能将原始空间映射为无穷维空间。如果 σ 选得很大的话, 高次特征上的权重实际上衰减得非常快, 所以近似一个低维的子空间; 反过来, 如果 σ 选得很小, 则可以将任意的数据映射为线性可分, 随之而来的可能是非常严重的过拟合问题。不过, 总的来说, 通过调控参数, 高斯核实际上具有相当高的灵活性, 也是使用最广泛的核函数之一。如下图:

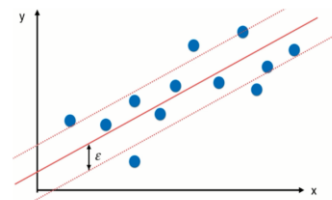
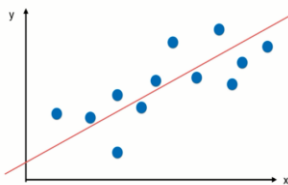


亲密度 affinity: 亲密度和距离正好相反。两点距离越远, 亲密度越低。常用高斯函数 $f(d) = \exp(-\gamma d^2)$ 。将距离 $\|x - y\|_p$ 转化成亲密度就得到了 RBF。从 $x - y$ 到 $\|x - y\|_p$ 到 RBF 是一个“向量 \rightarrow 距离(标量) \rightarrow 亲密度(标量)”的转化过程。

在 RBF 中其中 y (landmark) 的选取是任意的, 可以对每个数据样本点创建一个 landmark, 此时会将 $n \times D$ 的训练集变成 $n \times n$ 的训练集。另外地, 核函数 $\int K(x, y)\phi(x)dx = \lambda\phi(y)$, 这一点类似于特征向量的定义 $Ax = \lambda x$ 。

SVM 只能分类吗?

SVM 解决分类的时候, 同样指定一个 Margin 值, 只不过现在这个 Margin 的定义和解决分类问题中不同。我们想要 Margin 的两根虚线范围里能够包含的样本数据越多越好。这个思路其实和 SVM 解决分类问题是相反的思路。



文献阅读:

《统计学习方法》李航著;

吴恩达视频: <https://study.163.com/course/introduction/1210076550.htm>

5. 回归

线性回归:

直接求解

$y = h_{\theta}(X_i) = \theta^T X_i$, 式中 $\theta = (\theta_0, \theta_1, \dots, \theta_d)^T$, $X_i = (x_{i0}, x_{i1}, \dots, x_{id})^T$ (转置的原因是在概率论、机器学习等里面, 默认都是列向量), $1 \leq i \leq n$ 且 $i \in \mathbb{N}$ 。此式是输入特征的一个仿射变换 (affine transformation)。仿射变换的特点是通过加权和对特征进行线性变换 (linear transformation), 并通过偏置项来进行平移 (translation)。

其中 $x_0 \equiv 1$, 代表 θ_0 是偏置项。 i 代表的是对于第 i 行 (一个数据对象), d 代表的是有 d 个属性。因此 h 的作用是先是将数据集表示为一个 $n \times (d+1)$ 的矩阵 X (每一行代表一个示例, 后 d 个元素对应示例的 d 个特征值), 然后对矩阵的各个属性值赋予参数 θ , 求和得到一个目标预测值。

预测肯定存在误差, 即: 真实值=预测值+误差。对每个样本 $y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)}$ (标签=预测标签+误差项)。 $\varepsilon^{(i)}$ 独立同分布, 且服从 $\mu=0, \sigma=\sigma^2$ 的高斯分布 (正态分布)。由于误差服从

从高斯分布, 因此 $P(\varepsilon^{(i)}) = \frac{\exp\left(-\frac{(\varepsilon^{(i)} - \mu)^2}{2\sigma^2}\right)}{\sqrt{2\pi}\sigma} = \frac{\exp\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)}{\sqrt{2\pi}\sigma}$ (注意到均值 $\mu=0$)。因此 θ 与 x

组合, 成为真实值 y 的概率 $P(y^{(i)} | x^{(i)}; \theta) = \frac{\exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)}{\sqrt{2\pi}\sigma}$ 。将关注点从对 x 的概率函数转变为对 θ 的似然函数 (已知发生事件 x , 希望知道参数 θ 应该是多少)

$$L(\theta) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta),$$

$$\ln L(\theta) = \ln \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta) = \ln \prod_{i=1}^n \frac{\exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)}{\sqrt{2\pi}\sigma} = \sum_{i=1}^n \ln \frac{\exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)}{\sqrt{2\pi}\sigma} = n \ln \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2。$$

似然函数解决的是: 什么样的 θ 能使得这样的 θ 跟 x 组合之后, 成为 y 的可能性越大。

不妨令目标函数 $J(\theta) = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$, 因为 $\ln L(\theta)$ 越大越好, 因此 $J(\theta)$ 越小越好。

因为矩阵的平方等于转置乘自身, 因此 $J(\theta) = \frac{1}{2} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2 = \frac{1}{2} (X\theta - y)^T (X\theta - y)$, 求偏

$$\text{导 } \nabla_{\theta} J(\theta) = \nabla_{\theta} \left(\frac{1}{2} (\theta^T X^T - y^T) (X\theta - y) \right) = \nabla_{\theta} \left(\frac{1}{2} (\theta^T X^T X\theta - \theta^T X^T y - y^T X\theta + y^T y) \right) =$$

$$\frac{1}{2} (2X^T X\theta - X^T y - (y^T X)^T) = X^T X\theta - X^T y, \text{ 令偏导等于 } 0, \text{ 则 } \theta = (X^T X)^{-1} X^T y。$$

注 1: $P(x|\theta)$ 是一个有着两个变量的函数。如果将 θ 设为常量, 得到的是一个概率函数 (关于 x 的函数); 如果将 x 设为常量, 得到的是一个似然函数 (关于 θ 的函数)。已知一定量的样本数据 x , 在 x 条件下, 参数 θ 的最大似然估计就是使得该样本数据 x 出现的概率 (可能性) 最大, 那么参数 θ 的似然函数就是该样本数据出现的概率。似然和概率只是数值相等

$$L(\theta|x) = P(x|\theta) = \prod_{i=1}^n P(x_i|\theta), \text{ 但定义域不同, 似然函数的定义域是 } \theta, \text{ 而概率函数的定义域}$$

是 x , 即我们只能说“参数等于某个值时的似然是多少, 事件 (发生) 的概率是多少”。最大似然估计的含义: 事件已经发生, 里面未知参数取何值时, 该事件发生的概率最大。

注 2: $X^T X$ 是对称阵。对于对称阵 A , 有 $\nabla_{\theta}(\theta^T A \theta) = 2A\theta$ 。 $\nabla_{\theta} f(\theta)$ 即 $\frac{\partial f(\theta)}{\partial \theta}$ 。有时也对

目标函数除上常量 n 得到 $J(\theta) = \frac{1}{2n} \sum_{i=1}^n (\theta^T x_i - y_i)^2$, 不影响结果。

注 3: $J(\theta)$ 可以通过黑塞矩阵来说明存在极小值。因为 $\frac{\partial J}{\partial \theta_a} = \sum_{i=1}^n (\theta^T x_i - y_i) x_{ia}$,

$$\frac{\partial^2 J}{\partial \theta_a \partial \theta_b} = \sum_{i=1}^n x_{ia} x_{ib}, \text{ 因此黑塞矩阵是 } \begin{bmatrix} \sum_{i=1}^n x_{i1} x_{i1} & \cdots & \sum_{i=1}^n x_{i1} x_{in} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^n x_{in} x_{i1} & \cdots & \sum_{i=1}^n x_{in} x_{in} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n \end{bmatrix} \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} = X^T X, \text{ 因}$$

为 $x^T X^T X x = (Xx)^T (Xx) \geq 0$, 因此是半正定矩阵, 因此存在极小值。

注 4: $\text{MSE} = E(\theta^T) = (y - \theta^T X)^T (y - \theta^T X)$, $\frac{\partial \text{MSE}(\theta)}{\partial (\theta^T)} = 2X^T (X\theta^T - y)$ 亦得 $\theta^T = (X^T X)^{-1} X^T y$ 。

模型诊断 R2 分数: $R^2(y, \hat{y}) = 1 - \frac{SS_{res}}{SS_{tot}}$ 用于表示因变量(目标变量)的变异性(方差)能被自

变量(特征)解释的比例。式中, 残差平方和 $SS_{res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$, 总平方和 $SS_{tot} = \sum_{i=1}^n (y_i - \bar{y})^2$ 。

其中 y_i 为真实值, \bar{y} 为真实值的平均值, \hat{y}_i 为模型估计值。

梯度下降:

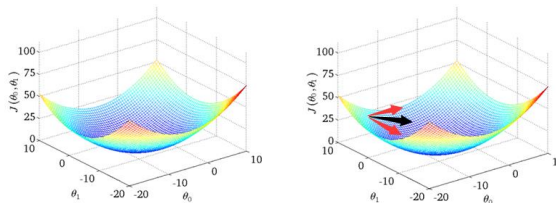
上面已经得到 θ 表达式, 那机器“学习”在哪里呢? 首先, 不是所有的算法都能直接计算出参数值。其次, 矩阵的逆不一定存在, 因此这个公式不能保证求出 θ 。

因为梯度的方向始终是向上, 而一般要求解的是损失函数的最小值, 所以叫梯度下降。

以只有 $\theta_{0,1}$ 参数为例, 因为 $x_{0,1}$ 没有关系, 那他们对应的 $\theta_{0,1}$ 更不可能有关系。因此

分别求偏导 $\frac{\partial J}{\partial \theta_0}, \frac{\partial J}{\partial \theta_1}$ (分别对应右图中的上/

下面的红色箭头)。向量相加即为综合结果。



目标函数: $J(\theta) = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)}))^2$ 。下面式中 i 代表第 i 个样本,

j 代表第 j 列(属性)。因为是对 θ_j 求导, 因此只需要考虑 θ_j 的系数 x_j 。

批量梯度下降 GD: $\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{n} \sum_{i=1}^n x_j^{(i)} (y^{(i)} - h_{\theta}(x^{(i)}))$, $\theta_j^{\cdot} = \theta_j + \frac{1}{n} \sum_{i=1}^n x_j^{(i)} (y^{(i)} - h_{\theta}(x^{(i)}))$ 。

随机梯度下降 SGD: $\theta_j^{\cdot} = \theta_j + (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$, 也就是 MBGD 中 batch_size=1。

小批量梯度下降 MBGD: $\theta_j^{\cdot} = \theta_j - \frac{\alpha}{n_c} \sum_{k \in n_c} x_j^{(k)} (h_{\theta}(x^{(k)}) - y^{(k)})$ 。 α 为学习率(learning rate),

一般越小越好。 n_c 为选择的数据个数(批量大小 batch_size), 常取 256, 512。这些可以调整但不在训练过程中更新的参数称为超参数(hyperparameter)。调参(hyperparameter tuning)是选择超参数的过程。超参数通常是我们根据训练迭代结果来调整的, 而训练迭代结果是在独立的验证数据集(validation dataset)上评估得到的。自动调参可参考: 自动机器学习 AutoML。

为什么学习率大了会在迭代过程中发散?

以 $J(\theta) = \theta^2$ 为例, $\frac{\partial J(\theta)}{\partial \theta} = 2\theta$, 迭代过程是 $\theta = \theta - \alpha \cdot 2\theta = (1 - 2\alpha)\theta$, 即 $\theta_k = (1 - 2\alpha)^k \theta_0$ 。

当 $\alpha \in (0, 0.5]$ 时, $0 \leq 1 - 2\alpha < 1$ (迭代不改变符号), $\lim_{n \rightarrow +\infty} \theta_n = 0$ (能收敛)。易证 $\alpha \in (0.5, 1)$ 时迭

代改变符号且能收敛。 $\alpha=1$ 时始终是 $\theta_{k+1}=-\theta_k$ 。 $\alpha>1$ 时迭代改变符号且不能收敛。

非线性回归(多项式回归):

一元 m 次多项式回归方程: $y = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_m x^m$

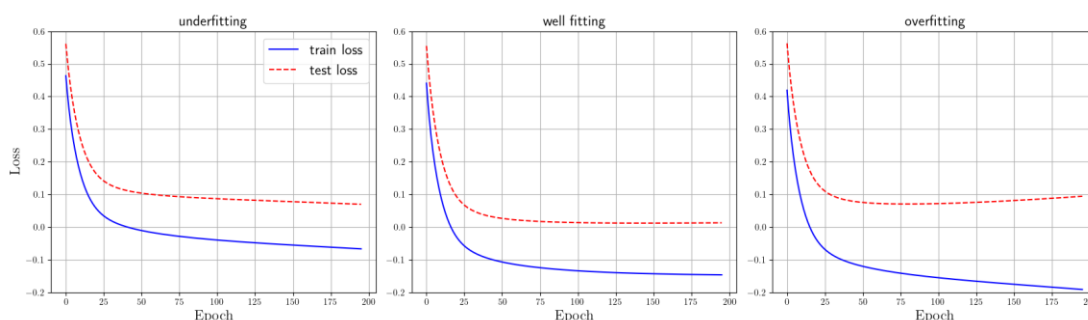
二元二次多项式回归方程: $y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2$

令 $x_1 = x, x_2 = x^2, \dots, x_m = x^m$ 可使一元 m 次方程化为 m 元线性回归方程 $y = \theta_0 + \theta_1 x_1 + \dots + \theta_m x_m$ 。因此用多元线性函数的回归方法就可解决多项式回归问题。

显然由于任意一个函数在一个较小的范围内都可以用多项式任意逼近,很容易出现过拟合问题。在线性回归中,如果训练集中数据的数量 N 小于等于数据特征数量 d ,就一定会出现过拟合现象,因为在 d 维空间中,任意小于等于 d 个点一定在同一个 $d-1$ 维超平面上。

处理过拟合问题常用的方法: 1.减少特征数量: 人工的挑选重要的特征,去除不重要的特征。2.正则化 Regularization: 保留所有特征,但是减少参数 ω 的值。

下图是模型训练和测试损失曲线示例,从左到右依次是欠拟合、恰好拟合和过拟合:



正则化: 避免过拟合,能处理特征点比样本点多的问题。用 α 代表惩罚力度:

L1 Lasso 套索回归 降维 有 $J = J_0 + \alpha \sum_{\omega} |\theta| = J_0 + \alpha \|\theta\|_1$, 解析解不便求。

$\frac{\partial J}{\partial w_i} = \frac{1}{N} \left(\frac{\partial J_0}{\partial w_i} + \alpha \text{sign}(w_i) \right)$ 在 $w_i \geq 0$ 时 $\frac{\partial J_0}{\partial w_i} + \alpha \geq 0 \Rightarrow J(w) \uparrow, w_i < 0$ 时 $\frac{\partial J_0}{\partial w_i} - \alpha < 0 \Rightarrow J(w) \downarrow$, 因此 $-\alpha \leq \frac{\partial J_0}{\partial w_i} \leq \alpha$ 时 $J(w)$ 在 w_i 处取得最小值。在逻辑回归中 $\frac{\partial J_0}{\partial w} = \sum_{i=1}^n (y_i - \sigma(xw)) x^{(i)}$, 代入得

$$\sum_{i=1}^n y_i x^{(i)} - \alpha \leq \sum_{i=1}^n \sigma(xw) x^{(i)} \leq \sum_{i=1}^n y_i x^{(i)} + \alpha, \text{ 式中 } 0 < \sigma(xw) < 1, \text{ 因此当 } \begin{cases} \sum_{i=1}^n y_i x^{(i)} \leq \alpha \\ \sum_{i=1}^n x^{(i)} - \sum_{i=1}^n y_i x^{(i)} \leq \alpha \end{cases}$$

时 $J(w)$ 在 $w=0$ 处取得最小值。简而言之,对于二值特征 $x_i \in \{0,1\}$,上式分别代表正/负样本中特征非零的个数。对于某个特征 w_j ,其正负样本中非零个数均 $\leq \alpha_j$ 时,该特征被忽略。

L2 Ridge 岭回归 采用惩罚算子 θ ,计算时用 $|y - X \cdot \theta|$, $\theta_{ridge} = \arg \min_{\theta} \left(\|y - X\theta\|_2^2 + \frac{\alpha}{2} \|\theta\|_2^2 \right) =$

$(X^T X + \alpha I)^{-1} X^T y$, 这是因为有 $J = J_0 + \alpha \sum \theta^2 = J_0 + \alpha \|\theta\|_2^2 = \frac{1}{2} (y - X\theta)^T (y - X\theta) + \frac{\alpha}{2} \|\theta\|_2^2$,

$$\frac{\partial J}{\partial \theta} = -X^T y + X^T X \theta + \alpha \theta = (X^T X + \alpha I) \theta - X^T y \Rightarrow \theta = (X^T X + \alpha I)^{-1} X^T y$$

为什么正则化能减少过拟合?

①以 $\theta_1 = [1 \ 0 \ 0 \ 0]$, $\theta_2 = \begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{bmatrix}$ 为例,特征数为4。使用L2回归时,对 θ_1 ,

惩罚项为 $\frac{1}{2} \cdot \alpha$ ；对 θ_2 ，惩罚项为 $\left[\left(\frac{1}{4}\right)^2 + \left(\frac{1}{4}\right)^2 + \left(\frac{1}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right] \cdot \alpha = \frac{1}{4} \alpha$ 。显然对 θ_1 惩罚值更大，

因此在 J_0 相同的情况下，能更倾向于选择参数 θ_2 。

换句话说，L2 衡量的是向量的大小，因此限制的是向量整体的大小，趋向于**均匀分布**。同样可得，L0 限制的是向量中非 0 元素的个数，趋向于**稀疏解**。不过 L0 含有示性函数 \mathbb{I} ，不可导，因此常用 L1 代替 L0。

②为什么①中的情况看起来是模型倾向于选择**更小数值** $\sum \theta^2$ 的参数组？这个可以反过来想，过拟合时次数更高的项的系数可能会变得更高。因为更高次数能更好地描述非线性变化，以此来拟合每个数据点，因此给其的权重更高，而正因其变化率太大的原因，系数也会变得更大。所以更小的系数 θ 可能意味着拟合得不够好，也就缓解了过拟合。

③过拟合的显著特征是微小的扰动带来结果的剧变，也就是**病态矩阵**。判断一个矩阵是不是病态矩阵，可以计算矩阵的**条件数**(矩阵的最大奇异值和最小奇异值之比)。如果矩阵存在很小的奇异值，那么它的逆就存在很大的奇异值^[证明 1]，这样对 y 中的微小变化会放大很多(可以用奇异值分解中最大的奇异值代表大部分信息来理解)。

$(X^T X + \alpha I)^{-1}$ 相对于 $(X^T X)^{-1}$ 是对所有奇异值加上 λ ^[证明 2]，因此避免了很小的奇异值。

[证明 1]: 对任意 $A: A = U \Sigma V^T \Rightarrow A^{-1} = V \Sigma^{-1} U^T$ ，也就是奇异值矩阵的各个元素取倒数。

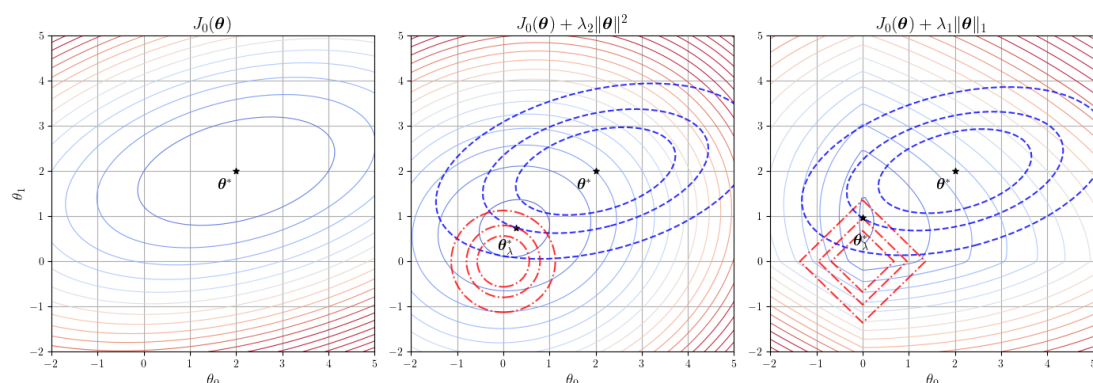
[证明 2]: $(X^T X)^{-1} = V \Sigma^{-1} U^T \xrightarrow{\alpha I = I @ (\alpha I) @ I^T} X^T X + \alpha I = V(\Sigma + \alpha I) U^T \Rightarrow (X^T X + \alpha I)^{-1} = V(\Sigma + \alpha I)^{-1} U^T$ 。

④L2 正则化损失函数相当于给原损失函数引入了一个关于 θ 的**高斯先验**。没看懂为什么会引入.....在 <https://zhuanlan.zhihu.com/p/361181741> 说在 PRML 第 1.2.5 节有，以后再看。

额外地，由于贯穿恒等式 $(\lambda I + PQ)^{-1} P = P(\lambda I + QP)^{-1}$, $P \in \mathbb{R}^{n \times m}$, $Q \in \mathbb{R}^{m \times n}$ ，因此 $f_\theta(X) = X(X^T X + \lambda I)^{-1} X^T y = XX^T (XX^T + \lambda I)^{-1} y = K(K + \lambda I)^{-1} y$ ，式中 $K = XX^T$, $K_{ij} = x_i^T x_j$ 是第 i, j 个样本的内积。因为样本本身所在的维度 \mathbb{R}^d 不一定能很好地反映样本本质的特征，所以需要特征映射函数 ϕ 把原始的特征经过某种筛选或者变换，我们称 $\kappa_{i,j} = \kappa(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ 为核矩阵 kernel matrix，称 $\kappa(\cdot, \cdot)$ 为核函数 kernel function。显然 $\phi(x) = x$ 就是我们这里得到的结果，也就是恒等映射。大多数情况下不需要显式地出现映射函数 ϕ ，直接用核函数 κ 进行计算即可。

不使用正则化，使用 L2，使用 L1 的对比如下图：

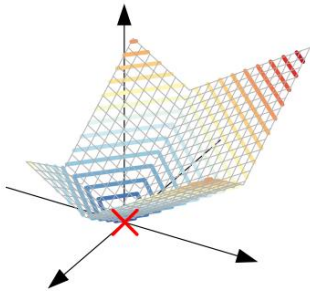
其中蓝虚线是 $J(\theta)$ ，红虚线是 $J(\theta) + \text{Regularization}$ ，彩线是整个损失函数。



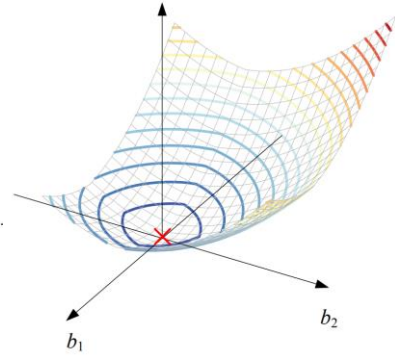
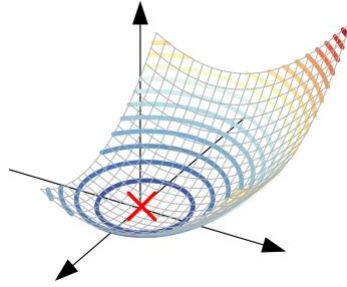
使用正则化后， θ 的解更靠近原点。而 L1 因为 0 处不可导，在 x, y 轴上存在明显的不平滑点，因此是趋向于稀疏解。

elastic net regression 弹性网络回归 以不同比例同时引入 L1 和 L2 正则项。正则化部分既有 L1 的尖点，也有 L2 的凸曲面。三种正则化方法的在特征数为 2 时的等高线曲面如下图：

L1 regularizer



L2 regularizer



Logistic 回归:

以购房为例，假如只考虑能不能买得起房，就是一个突变问题，需要 S 型曲线来拟合。

即 logistic 函数/Sigmoid 函数 $\sigma(z) = \frac{1}{1+e^{-z}} = \frac{e^z}{1+e^z}$ ，值域为 (0,1)，可用来代表概率。虽然名为回归，但实际常解决二分类问题。一般用 σ 表示 softmax/logistic 函数。

将线性回归中的 $\theta^T X$ 代入公式得 $h_\theta(X) = \sigma(\theta^T X) = \frac{1}{1+e^{-\theta^T X}}$ ，则 $\begin{cases} P(y=1|x;\theta) = h_\theta(x) \\ P(y=0|x;\theta) = 1-h_\theta(x) \end{cases}$ ，

整合得 $P(y|x;\theta) = (h_\theta(x))^y (1-h_\theta(x))^{1-y}$ (代入 $y=0,1$ 即可验证此公式正确性)。依旧取似然函数

$L(\theta) = \prod_{i=1}^n (h_\theta(x_i))^{y_i} (1-h_\theta(x_i))^{1-y_i}$ ，取对数 $\ln L(\theta) = \sum_{i=1}^n (y_i \ln h_\theta(x_i) + (1-y_i) \ln(1-h_\theta(x_i)))$ ，为了符合

梯度下降，因此这里取负。引入目标函数 $J(\theta) = -\frac{1}{n} \ln L(\theta)$ ，对其求偏导 $\frac{\partial J(\theta)}{\partial \theta_j}$

$$= -\frac{1}{n} \sum_{i=1}^n \left(y_i \frac{1}{h_\theta(x_i)} \frac{\partial h_\theta(x_i)}{\partial \theta_j} + (1-y_i) \frac{-1}{1-h_\theta(x_i)} \frac{\partial h_\theta(x_i)}{\partial \theta_j} \right) = -\frac{1}{n} \sum_{i=1}^n \left(y_i \frac{1}{h_\theta(x_i)} + (1-y_i) \frac{-1}{1-h_\theta(x_i)} \right) \frac{\partial h_\theta(x_i)}{\partial \theta_j}$$

$$= -\frac{1}{n} \sum_{i=1}^n \left(y_i \frac{1}{\sigma(\theta^T x)} + (1-y_i) \frac{-1}{1-\sigma(\theta^T x)} \right) \frac{\partial \sigma(\theta^T x)}{\partial \theta_j}$$

$$= -\frac{1}{n} \sum_{i=1}^n \left(y_i \frac{1}{\sigma(\theta^T x)} + (1-y_i) \frac{-1}{1-\sigma(\theta^T x)} \right) \left(\sigma(\theta^T x)(1-\sigma(\theta^T x)) \frac{\partial \theta^T x_i}{\partial \theta_j} \right)$$

$$= -\frac{1}{n} \sum_{i=1}^n (y_i(1-\sigma(\theta^T x)) - (1-y_i)(\sigma(\theta^T x))) x_j^i = -\frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\theta^T x)) x_j^i，因此参数更新为：$$

$$\theta_j = \theta_j + \alpha \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\theta^T x)) x_j^i = \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1+e^{-\theta^T x}} - y_i \right) x_j^i。$$

为什么逻辑回归使用 sigmoid 作为预测结果的输出？

逻辑回归的输出服从伯努利分布 $P(y|x;w) \sim B(\phi)$ ，可以写成指数族分布的形式：

$$P(y|x;w) = P(y;\phi) = \phi^y (1-\phi)^{1-y} = \exp(y \log \phi + (1-y) \log(1-\phi)) = \exp\left(y \log \frac{\phi}{1-\phi} + \log(1-\phi)\right)$$

指数族分布 $P(y;\eta) = b(y) \exp(\eta^T T(y) - a(\eta))$ ，式中 $b(y)$ 是基础度量值， η 是自然参数 (可以使向量也可以是标量)， $T(y)$ 是充分统计量， $\exp(-a(\eta))$ 起到归一化作用。正态、伯努利、多项式分布都是指数族分布。指数族分布的对数似然函数通常是凸函数。

将伯努利分布的指数族形式代入指数族分布中可得：

$$b(y) = 1, \eta = \log \frac{\phi}{1-\phi}, a(\eta) = -\log(1-\phi) = \log(1+e^\eta), b(y) = 1$$

$$\text{因此 } \phi = \frac{1}{1+e^{-\eta}}, \text{ 因此 } P(y=1|x;w) = \phi = \frac{1}{1+e^{-\eta}}。$$

对于多分类问题，使用 softmax: $\hat{y}_i = \frac{e^{(X \cdot w_i + b_i)}}{\sum_{k=1}^C e^{(X \cdot w_k + b_k)}} = \frac{e^{(X \cdot w_i)}}{\sum_{k=1}^C e^{(X \cdot w_k)}}$ 。

当 $C=2$ 时 softmax 就退化为 logistic 函数。

注：上文中有一些失误，一般而言使用的是 Xw 而不是 wX ，因为 $X_{n \times D} w_{D \times m} = y_{n \times m}$ 这个形状才是符合需求的，这样才将输入的 n 个特征变为输出的 m 个特征。

损失函数(交叉熵): $J(\theta) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^C y_k^{(i)} \ln \hat{y}_k^{(i)}$ ， $y_k^{(i)}$ 是第 i 个样本属于第 k 类的指示函数。

3. softmax是对上面介绍的映射的误称（虽然深度学习领域中很多人都使用这个名字）。真正的softmax被定义为 $\text{RealSoftMax}(a, b) = \log(\exp(a) + \exp(b))$ 。

1. 证明 $\text{RealSoftMax}(a, b) > \max(a, b)$ 。
2. 证明 $\lambda^{-1} \text{RealSoftMax}(\lambda a, \lambda b) > \max(a, b)$ 成立，前提是 $\lambda > 0$ 。
3. 证明对于 $\lambda \rightarrow \infty$ ，有 $\lambda^{-1} \text{RealSoftMax}(\lambda a, \lambda b) \rightarrow \max(a, b)$ 。
4. soft-min会是什么样子？
5. 将其扩展到两个以上的数字。

6. 关联规则挖掘

例子如右：这里研究关联规则：“啤酒→尿布”

置信度 confidence: $C(\text{啤酒} \rightarrow \text{尿布}) = P(\text{尿布} | \text{啤酒}) = P(\text{啤酒} \cup \text{尿布}) / P(\text{啤酒}) = \frac{3}{4} = 75\%$

支持度 support: $S(\text{啤酒} \cup \text{尿布}) = N(\text{尿布与啤酒同时出现}) / N(\text{全部购物数据}) = \frac{3}{100}$

提升度 lift: $L(\text{啤酒} \rightarrow \text{尿布}) = \text{confidence}(\text{啤酒} \rightarrow \text{尿布}) / \text{support}(\text{尿布}) = 25$

$$\text{即 } C(A \rightarrow B) = P(B|A) = \frac{P(AB)}{P(A)}, S(AB) = \frac{N(AB)}{N} = P(AB), L(A \rightarrow B) = \frac{C(A \rightarrow B)}{S(B)} = \frac{P(AB)}{P(A)P(B)}.$$

为什么需要提升度？如下例，数据如右表：

对于 $C(\text{舟} \rightarrow \text{原}) = P(\text{都玩}) / P(\text{舟}) = 0.8$ ，说明玩方舟的大部分玩原。

对于 $C(\text{不舟} \rightarrow \text{原}) = P(\text{不舟但原}) / P(\text{不舟}) = 0.7$ ，说明不玩舟的大部分玩原。

那玩不玩原，就和玩舟的关系不大。此时提升值 $L(\text{舟} \rightarrow \text{原}) =$

$C(\text{舟} \rightarrow \text{原}) / P(\text{原}) = 0.8 / 0.75 = 1.0667$ ，此值并不大。（若 $Lift=1$ ，则代表二者独立）

简而言之，置信度= C 表明在前件发生时会有 C 的概率会发生后件，支持度= S 表明规则的覆盖范围是 S （也就是一起出现的概率），提升度= L 表明在前件发生时后件发生的可能性是其普通情况下的 L 倍。

设置**最小支持度阈值**和**最小置信度阈值**，当某个关联规则同时满足这两个预选设好的阈值，便把它称为**强关联规则**。

	玩原神	不玩原神	
玩方舟	80	20	100
不玩方舟	70	30	100
	150	50	200

事务: 每一条购物信息 **项:** 购物信息中的一项物品 **项集:** 包含 0 或多个项的集合

k 项集: 含有 k 个项的项集 **关联规则:** 蕴含式 $A \rightarrow B$ ，其中 $A \neq \emptyset, B \neq \emptyset, A \cap B = \emptyset$

频繁项集: 支持度大于预定义的最小支持度阈值的项集

于是，从数据集中挖掘强关联规则的过程分为两步：

1. 需要找出满足最小支持度阈值的项集，即频繁项集；
2. 根据最小置信度阈值，从频繁项集中生成强关联规则。

也就是找出 1 项集，2 项集，...， n 项集。然后从中找出大于支持度阈值（如 50%）的项集，再判断置信度。

当数据集的项过多时，该算法产生的**候选项集**过多，计算量太大。假设有 100 个 1 项集，

则候选频繁项集个数： $C_{100}^1 + C_{100}^2 + \dots + C_{100}^{100} = 2^{100} - 1 = 1267650600228229401496703205375$

此时要用到 **Apriori 性质**：

如果一个项集 A 是频繁项集，那么它的非空子集 B 也是频繁项集。

或者：如果一个项集不是频繁项集，那它的超集也不是频繁项集。

因此只需检索由频繁 k 项集生成的 $k+1$ 项集，而若某 k 项集已不满足， $k+1$ 项集就更不可能。

缺点:

Apriori 算法每一次迭代都需要遍历数据集，当数据集过大无法一次性载入内存时，遍历的效率会变得非常低。Apriori 算法仍然会产生大量无用的候选多项集。

FP-growth 算法:

高效地发现频繁项集，但是不能用于发现关联规则。FP-growth 算法的执行速度快于 Apriori 算法，通常性能要好两个数量级以上。FP-growth 算法只需要对数据集扫描两次，它发现频繁项集的过程如下：构建 FP 树，从 FP 树中挖掘频繁项集。

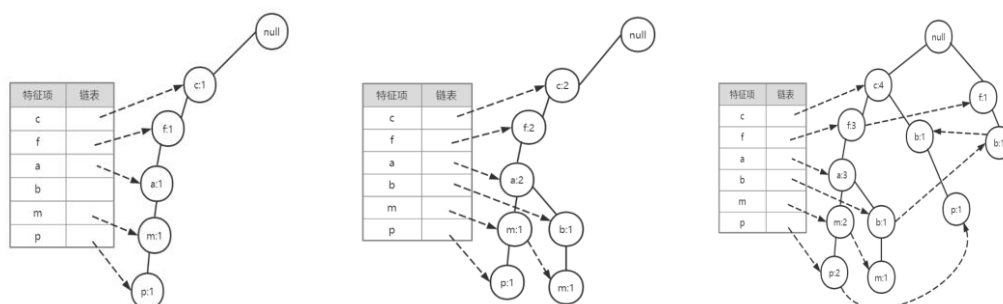
如下例，原数据为前两列：

编号	含有数据项	排序的频繁项集
1	f,a,c,d,g,i,m,p	c,f,a,m,p
2	a,b,c,f,l,m,o	c,f,a,b,m
3	b,f,h,j,o	f,b
4	b,c,k,s,p	c,b,p
5	a,f,c,e,l,p,m,n	c,f,a,m,p

特征项	链表
c	
f	
a	
b	
m	
p	

null

- ①设置最小支持度为 0.6，频繁 1 项集应出现至少 $0.6 * 5 = 3$ 次。
- ②扫描数据库(上表的前两列)，识别所有的单个项并计算它们的支持度计数，提取出支持度计数 ≥ 3 的项，并按照支持度从高到低排序： $L = \{\{c:4\}, \{f:4\}, \{a:3\}, \{b:3\}, \{m:3\}, \{p:3\}\}$ ，得到上表的第三列。
- ③构建根节点，标记为 null，创建频繁项集，将链接列设为空，如右上图。
- ④将第 1 条事务插入频繁模式树中，只插入频繁项集部分 $\{c,f,a,m,p\}$ ，将每个节点的支持度计数设为 1，同时设置好对应项的链接。如左下图。
- ⑤将第 2 条事务插入频繁模式树中 $\{c,f,a,b,m\}$ ，当节点相同，共享节点并将节点的支持度加 1，同时更新对应项的链接。如下中图。
- ⑥最终频繁模式树如右下。



⑦条件模式基(conditional pattern base)是指从原始数据集的频繁树中投影出的子数据集，它与特定后缀相关联。以所查询元素项为结尾的路径集合。每条路径其实都是一条前缀路径(prefix path)。总之，一条前缀路径是介于所查询元素项与树根节点之间的所有内容。

通过创建的头指针表来得到前缀路径。头指针表包含相同类型元素链表的起始指针。一旦到达了每一个元素项，就可以上溯这棵树直到根节点为止。

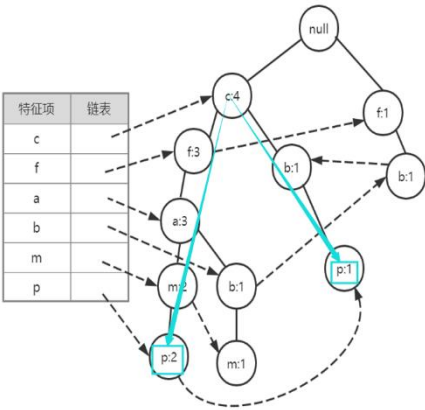
条件模式基有什么用？以 p 的条件模式基为例，当 p 元素出现的前提下，c,f,a,b,m 各元素出现的概率。

例如：构造以 p 为后缀的条件模式库的方法如下：

- 1.找到频繁模式树中所有的 p 节点，如图蓝框。
 - 2.回溯所有的路径(如图)得到(cbp, cfamp)，除去 p 节点后，将路径中的剩余节点当作新的项集。
 - 3.项集的支持度设置为对应路径中 p 节点的支持度。
- 得到以 p 为后缀的条件模式库为： $\{cfam: 2, cb: 1\}$ 。寻找方法如下图：

并以此可得到其他特征项的条件模式库：

特征项	条件模式库
c	{}
f	{c:3}
a	{cf:3}
b	{cfa:1; c:1; f:1}
m	{cfa:2; cfab:1}
p	{cfam:2; cb:1}



7. 推荐系统

矩阵分解 MF(matrix factorization):

假设 N 个用户, M 部电影有 D 个特征。我们预测 $R = \underset{N \times M}{P} \underset{N \times D}{Q} \underset{D \times M}$, 式中 R 是预测 N 个用户对 M 部电影的评分, P 是用户喜好矩阵, Q 是电影特征矩阵。

因为实际获取到的一般是稀疏矩阵 R , 我们需要反推 P, Q 。因此损失函数就是预测的 R' 与真实的 R 之间的误差, 可选 MSE。注意正则化的时候需要考虑每个用户对电影打分的数量不同, 因此损失函数是 $J(P, Q) = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^M \mathbb{I}_{ij} \left((r'_{ij} - r_{ij})^2 + \lambda (\|p_i\|^2 + \|q_j\|^2) \right)$, λ 受示性函数 \mathbb{I} 控制, \mathbb{I}_{ij} 代表用户 i 是否对电影 j 评分。

用户 i 对电影 j 的评分 r_{ij} 是用户偏好与电影特征的内积 $r_{ij} = p_i q_j$ 。向量内积是**双线性函数 bilinear model**。双线性的含义为, 二元函数固定任意一个自变量时, 函数关于另一个自变量线性。

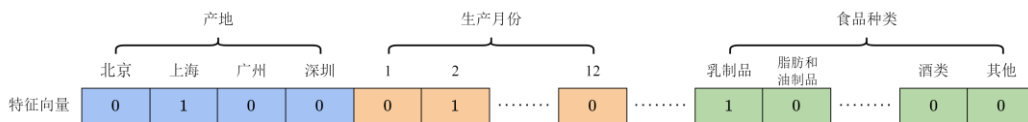
因子分解机 FM(factorization machines)

MF 的目标是从交互的结果中计算出用户和物品的特征。而 FM 则正好相反, 希望通过物品的特征和某个用户点击这些物品的历史记录, 预测该用户点击其他物品的概率, 即**点击率(click through rate, CTR)**。

因为物品之间可能存在关联, 所以也要引入双线性模型: $y(x) = \sum_{i=0}^D \theta_i x_i + \sum_{i=1}^{D-1} \sum_{j=i+1}^D w_{ij} x_i x_j$,

即 $y(x) = \theta^T x + \frac{1}{2} x^T W x$, 其中 $x_0 \equiv 1$ 对应偏置项。

如果将物品离散编码, 可以使用**多域独热编码 multi-field one-hot encoding**, 将物品的多个特征各自独热编码成的向量依次拼接起来。比如:



在此编码方式下, $\frac{\partial y}{\partial w_{ij}} = x_i x_j$ 在大多数时候等于 0, 无法更新参数。矩阵分解 $W = VV^T$ 后

能解决这一问题, 其中 $W \in \mathbb{R}^{D \times D}, V \in \mathbb{R}^{D \times k}$ 。

算法优化后时间复杂度 $O(kd)$: $y(x) = \sum_{i=0}^D \theta_i x_i + \frac{1}{2} \sum_{l=1}^k \left(\left(\sum_{i=1}^D v_{il} x_i \right)^2 - \sum_{i=1}^D v_{il}^2 x_i^2 \right)$, 化简过程见

<https://hml.boyuai.com/books/chapter7> 第 7.3 节。

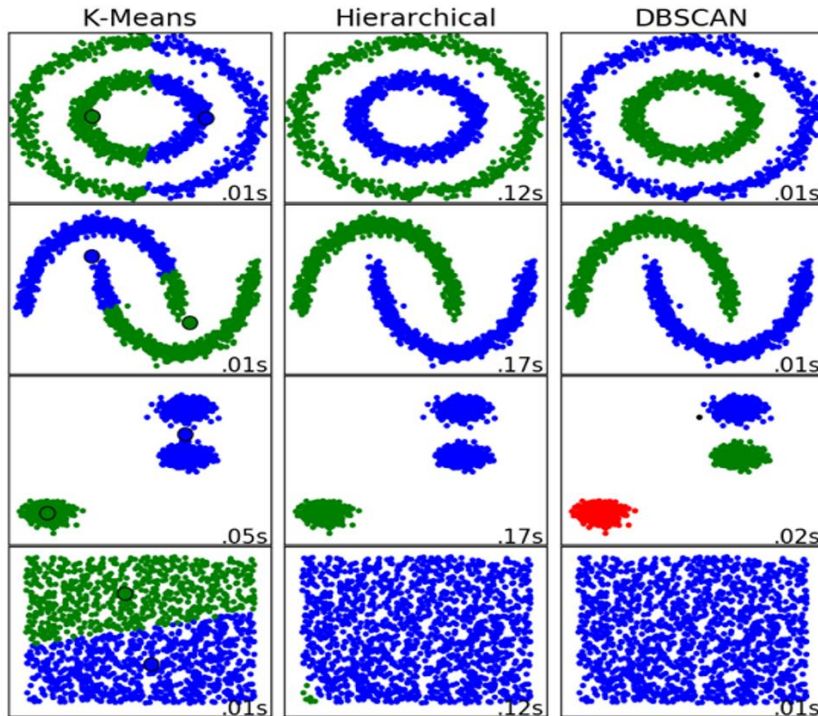
聚类分析

0.说在前头

聚类模型的**本质**：将数据集中相似的样本(无标签)进行分组(加标签)的过程。

每个组称为一个簇(cluster)。每个簇的样本对应一个潜在的类别。样本是没有类别标签的，因此聚类是一种典型的无监督学习任务，这些簇满足以下两个条件：相同簇的样本之间距离较近，不同簇的样本之间距离较远。

常见聚类：划分聚类、层次聚类、基于密度。



1.划分聚类 (基于原型的聚类 Prototype-Based Clustering)

假设聚类结构能通过一组原型刻画。

K 均值聚类：

①**取初始质心**：从数据集中随机选取 k 个初始质心。(k 是事先给定的)

②**划分数据点**：计算每个点与各个质心的距离(如欧氏距离)，将数据点划分到距离最近的质心所在的簇内。

③**寻找新质心**：原质心不在新簇的质心位置，于是应该计算每一个簇的新质心，也就是求均方误差和 SSE(Sum of Squared Error)最小化问题。 $SSE = \sum_{i=1}^k \sum_{x_j \in C_i} r_{ik} \|x_j - u_i\|^2$ ，使 SSE 最小

的中心点是均值 $u_i = \frac{\sum_{x_j \in C_i} x_j}{j_i}$ 。式中 C_i 是第 i 个簇，含有 j 个数据点， x_j 为该簇内第 j 个数据

点， u_i 是质心， $r_{ik} \in \{0,1\}$ (若 x_i 被划分到簇 k 中则 $r_{ik}=1$ ，且对于 $j \neq k$ 有 $r_{ij}=0$ ， $\sum_{j=1}^k r_{ij}=1$)。 $\sum_{i=1}^k$

表示各个簇， $\sum_{x_j \in C_i}$ 表示簇内。

④重复②~③直至各个数据点不再变更所属的簇。

多次运行，每次选取一组不同的随机初始质心可以避免陷入局部最优。

评价标准：

1. inertia_ – Sum of squared distances of samples to their closest 样本离最近聚类中心的总和。

$$2. \text{轮廓系数 } S_i = \frac{b_i - a_i}{\max\{a_i, b_i\}} = \begin{cases} 1 - \frac{a_i}{b_i}, & a_i < b_i \\ 0 & \\ \frac{b_i}{a_i} - 1, & a_i > b_i \end{cases}, \text{ 式中 } a_i \text{ 为样本 } i \text{ 到同簇其他样本的平均距离, } b_{ij}$$

为样本 i 到某簇 j 所有样本的平均距离, $b_i = \min\{b_{i1}, b_{i2}, \dots, b_{ik}\}$ 。理想情况下 a_i 很小, b_i 很大。

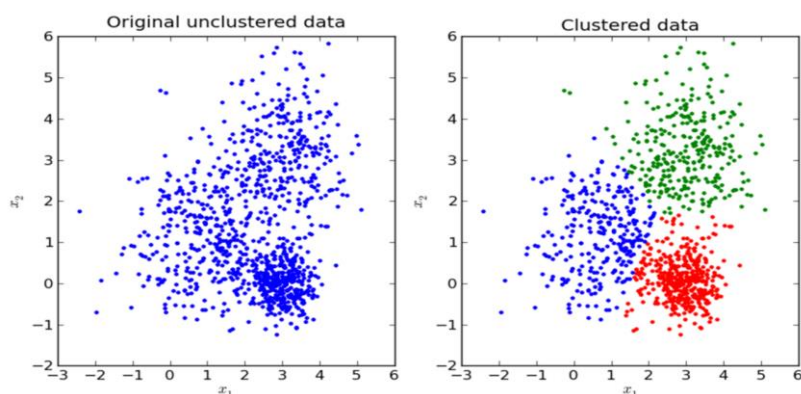
因此 S_i 接近 1 说明样本 i 聚类合理, S_i 接近 0 说明在两簇的边界上, S_i 接近 -1 代表样本 i 更应该分类到另外的簇中。

K-means 算法的优点：算法实现简单、直观。

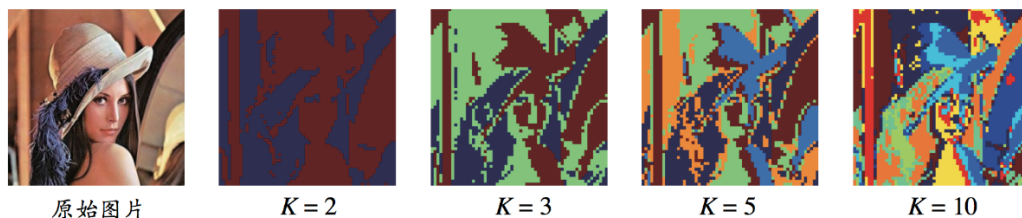
K-means 算法的缺点：事先指定 k 值，聚类结果依赖于 k 个初始质心的选择。易陷入局部最优，不易处理非簇状数据。聚类结果容易受离群值影响。

时间复杂度： $O(mnkt)$ ： m 为特征维度、 n 为样本数量、 k 为簇数量、 t 为迭代次数。每一步迭代中，每个样本需要与 k 个簇中心进行距离计算。

K 均值聚类分类效果如下：



图像压缩：给出了原始图像以及使用不同 k 值得到的 K-means 分割结果：



其他原型聚类

K-means 聚类变种——K-medoids 聚类（K 中心点聚类）

K-medoids 聚类使用数据点作为聚类中心，而不是像 K-means 聚类那样使用平均值。这意味着 K-medoids 对于噪声和异常值更具鲁棒性，并且能够产生更好的聚类结果。

鲁棒性是指一个算法或系统对于异常情况、噪声或错误的输入数据的处理能力和稳定性。

具有良好鲁棒性的算法或系统能够在面对这些问题时仍然保持稳定可靠的输出结果。

二分 K-means: 一种基于 K-means 算法的聚类方法, 与传统的 K-means 不同之处在于, 它利用了二分思想, 将原始数据集**不断划分为两个子集**, 直至达到指定的聚类数目, 得到最终的聚类结果。(每次从所有簇中选择具有最大 J 值的簇, 然后划分为两个簇, 重复该过程, 直至簇集合中含有 K 个簇。其中 J 指的是簇内平方误差和 SSE, K 指的是最终聚类的簇数)。

K-medians 聚类: 使用样本每个维度的中位数作为簇的质心。

K-means++ 聚类: 在初始化质心的过程中, 选择相互距离尽可能远的样本作为初始质心。

基于粗糙集的 K-means 聚类: 在基于粗糙集的 K-means 算法中, 一个样本可以被划分给多个簇。

学习向量量化 LVQ(Learning Vector Quantization): 是一种有监督的聚类方法, 因为它需要训练数据的标签信息(或假设存在数据带有类别标签)来指导向量更新过程。它通常用于具有清晰类别的数据集。可看作通过聚类来形成类别“子类”结构, 每个子类对应一个聚类簇。

高斯混合聚类 GMM(Gaussian Mixture Model): 是一种无监督学习方法, 使用高斯分布、贝叶斯公式、极大似然法、聚类等原理, 用于探索数据中的潜在聚类结构。

2. 层次聚类

层次聚类(hierarchical clustering)在不同层级上对样本进行聚类, 逐步形成树状的结构。根据层次分解是自底向上(合并)是自顶向下(分裂)方式, 层次聚类方法可以分为: **聚合式(凝聚式)聚类(agglomerative clustering)**、**分拆式(分裂式)聚类(divisive clustering)**, 两种方法均是启发式的策略, 并没有去优化一个明确的目标函数来实现聚类, 很难严格评价聚类的效果。

层次聚类方法的**局限性**: 1. 计算量大。2. 一旦合并或分裂执行, 就不能修正。如果某个合并或分裂决策在后来证明是不好的选择, 该方法也无法退回并更正。

原型聚类样本所在簇可能变化, 但层次聚类样本所在簇迭代后不变。

凝聚层次聚类

采用自底向上策略, 首先将每个样本作为单独的一个原子簇, 然后合并这些原子簇形成越来越大的原子簇, 直到所有的样本都在一个簇中(层次的最上层), 或者达到一个终止条件。绝大多数层次聚类方法属于这一类。

在开始时把每个样本都当成一簇, 然后在每一次迭代中将最相似的(距离最近)两个簇进行合并, 直到把所有簇合并为包含所有样本的一簇。

单连接(最近邻距离): 簇 G 和簇 H 之间的距离定义为两簇之间最近的成员之间的距离。

完整连接(最远邻距离): 簇 G 和簇 H 之间的距离定义为两簇之间最远的成员之间的距离。

平均连接: 表示两簇间所有成员对的平均距离 $d_{avg}(G, H) = \frac{\sum_{i \in G} \sum_{i' \in H} d_{i,i'}}{n_G n_H}$

单连接法只需要两簇内有成员对距离足够近就将两簇合并, 而并没有考虑其他簇内其他成员的距离, 因此单连接方法形成的簇很有可能违背**紧致性特征**(簇内成员应该尽可能相似)。(也就是因为存在某些点相似就把两个簇当做相似了)

完整连接法是另外一个极端, 只有当两簇的联合的成员间的距离相对较小时才将两簇进行合并, 因此完整连接法倾向于生成紧致簇。(与单连接恰好相反)

平均连接法介于单连接和完整连接之间, 易于生成相对紧致的簇同时簇间距离较远。

AGNES 层次聚类算法(单连接自底而上凝聚层次聚类方法):

先将每个样本看成一个簇, 然后每次对距离最短的两个簇进行合并, 不断重复, 直到达

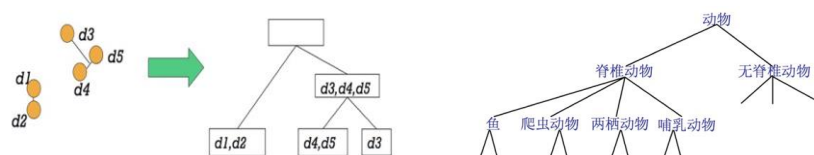
到预设的聚类簇个数。

分裂层次聚类

采用自顶向下策略, 首先将所有样本置于一个簇中, 然后逐渐细分为越来越小的簇, 直到每个样本自成一个簇, 或者达到某个终止条件, 例如达到了某个希望的簇的数目, 或者两个最近的簇之间的距离超过了某个阈值。

分拆式聚类将所有样本集合看作一簇, 以自上而下的方式, 递归地将现有的簇分拆为两个子簇。如: **二分 K-means 聚类**(选择半径最大的簇, 对该簇进行 K-means 聚类分为两个子簇, 重复此过程直到达到想要的簇个数)。**最小生成树法**(将每个样本看作一个图节点, 将样本间距离看作节点间边的权重, 根据此图建立最小生成树。从权重最大处将该簇分拆为两簇, 然后重复此过程直到达到想要的簇个数。实际上, 该方法得到的聚类结果和单连接的聚合聚类得到的结果一致。)

聚合式聚类, 分拆式聚类如下图(左聚合 右分拆):



3. 密度聚类

密度聚类又称为基于密度聚类(density-based clustering)。此类算法假设聚类结构由样本分布的紧密程度确定, 以数据集在空间分布上的稠密程度为依据进行聚类, 即只要一个区域中的**样本密度大于某个阈值**, 就把它划入与之相近的簇中。

不同于基于质心的原型聚类算法, 基于密度的算法不需要事先设置 k 值。常用的密度聚类算法: DBSCAN、MDCA、OPTICS、DENCLUE 等, 下面主要讲解 DBSCAN 算法。

① 用户设定邻域参数(ϵ , Minpts)。

ϵ 为半径, Minpts 代表一个簇中最少的数据点个数, 高于该值为高密度区域, 否则为低密度区域。

② 数据点分类

基于邻域参数, 在 DBSCAN 算法中, 所有的数据点可以分为以下三个类别:

核心点: 如果某个点的邻域内的数据点数目高于阈值 Minpts, 则将这个点视为核心点。

边界点: 边界点是位于核心点邻域之内的, 但是其自身的邻域内数据点数小于 Minpts 的点, 起到将高密度区域与低密度区域分割开的作用。

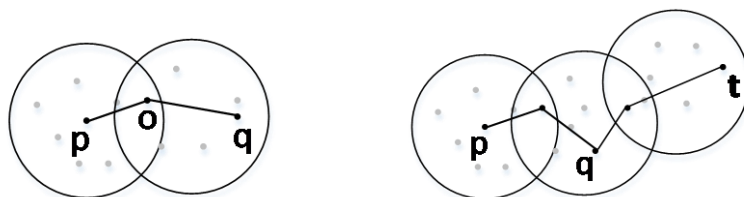
噪声点: 既不是核心点也不是边界点的其他数据点, 他们组成低密度区域。

密度直达: o 在 p 的邻域内, 且 p 为核心点, 则称从 p 到 o 是密度直达。

密度可达: o 在 p 的邻域内, 从 p 到 o 是密度直达, 而 q 对象的邻域内不包括 p , 但是包括 o , 这样 $p \rightarrow o \rightarrow q$, 称 p 到 q 是密度可达的。

密度相连: q 和 p 是密度可达的, q 和 t 也是密度可达的, 则 p 和 t 是密度相连的。

下图分别表示 p 到 q 是密度可达, p 和 t 是密度相连。



DBSCAN 的目标为找到密度相连数据点的最大集合, 此集合作为最终的一簇。其实就是从某个随机点开始向外扩张直至 ϵ 邻域内的数据个数小于 Minpts。

首先核心点各自成簇，采用密度相连的概念逐步对簇进行合并，打上标记。最终核心点密集的区域会被低密度噪声点包围，噪声点不单独成簇。所以采用 DBSCAN 进行聚类，最终的结果有一些数据点是没有标记的(作为簇-1)，这些就是噪声点。

优点：

1. 无需事先设定簇的个数，算法根据数据自身找出各簇。
2. 适于稠密的非凸数据集，可以发现任意形状的簇。
3. 可以在聚类时发现噪音点、对数据集中的异常点不敏感。
4. 对样本输入顺序不敏感

缺点：

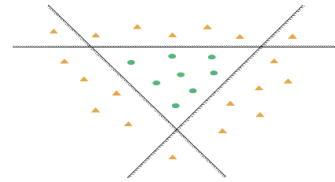
1. 因为是基于密度，若两个簇没有明显的可分间隔，很有可能被合并为同一个簇。
2. 参数调节较为复杂，参数设置对结果影响较大。
3. 对高维的数据处理效果不好。

集成模型

1.定义与作用与一些术语

集成学习 Ensemble Learning 通过构建并结合多个学习器来完成学习任务。

作用：①增强模型的表达能力：单个感知机无法正确分类数据，集成三个感知机能正确分类数据，如右图。②降低误差：分类器之间独立。



如果所有的个体学习器都是同类型的，则称这样的集成是**同质的**，同质集成的个体学习器称为基学习器。如果个体学习器不是同类型的，那么这样的集成是**异质的**，异质集成的学习器称为组件学习器。

弱学习器：性能略优于随机猜测的学习器。如在二分类问题上预测精度略高于 50% 的分类器。

2.如何得到结果

①**多数投票方法**(majority vote)：选取预测次数最多的类别 C_i (有时拒绝预测没有过半的类别)。也能使用加权投票。

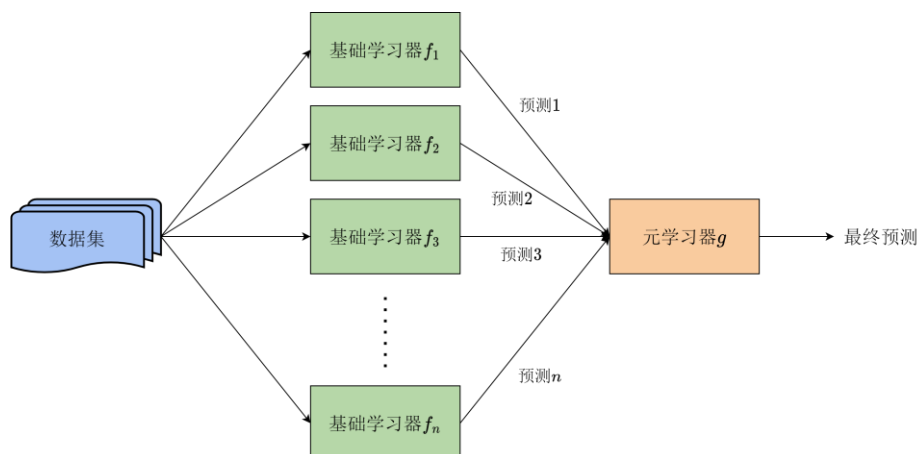
②**平均**(averaging)： $H(x) = \frac{1}{n} \sum_{i=1}^n h_i(x)$

加权平均(weighted averaging)： $H(x) = \sum_{i=1}^n w_i h_i(x)$ ，其中 $\sum_{i=1}^n w_i = 1$ 且 $w_i \geq 0$ 。如 AdaBoost。

个体学习器性能差异大的时候选用加权平均，个体学习器差异小的时候选用算术平均。因为权重是训练得出的，不完全可靠。

③**学习法**：通过另一个学习器来结合预测结果。比如 stacking 将弱学习器的学习结果作为输入，重新训练一个学习器来得到最终结果，即：多个初级学习器(基本模型)被训练，它们的输出作为次级学习器(元模型)的输入，然后次级学习器被训练。

3.典型的集成学习器 ensemble learner

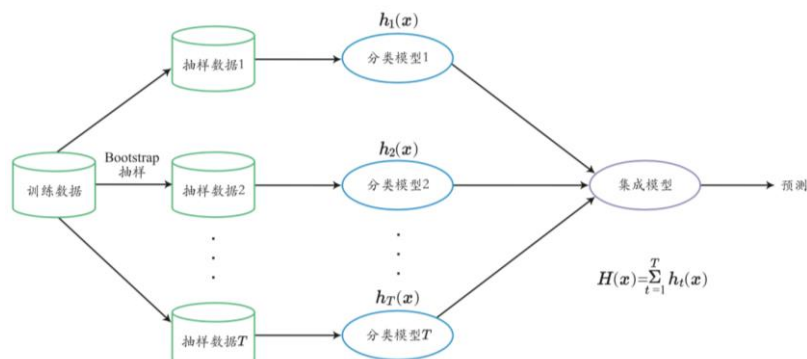


通常将底层的模型 f 称为基础学习器 base learner，而将称为元学习器 meta learner。

Bagging：

自举聚合(bootstrap aggregation, bagging)算法, 其中自举是自举采样(bootstrap sampling)。

对样本或特征随机取样(自举), 学习产生多个独立的基础模型(注意需要是同一种类的模型), 然后平均所有模型的预测值。基于每个抽样数据集 D_i 训练得到基础模型 $h_i(x)$, 最终的分分类器 $H(x) = \text{sign}\left(\sum_{i=1}^T h_i(x)\right)$, 如下图。典型代表随机森林。



自举采样 Bootstrap: 为了保证随机性, 尽可能降低不同子数据集之间的相关性, 采用了允许重复的有放回采样。通过有放回抽样 n 次得到与原数据集 D 大小一样的样本集

D_i 。一个样本有 $1 - \frac{1}{n}$ 的概率不被选到, 则一个样本不被抽到的概率为

$(1 - \frac{1}{n})^n \approx \frac{1}{e} = 0.368$, 说明训练集包含 63.2% 的不同样本。每个 D_i 学习得到一个分类器

M_i 。对于分类问题采用多数投票, 对于回归问题采用均值投票, 输出集成模型:

$$H(x) = \begin{cases} \frac{1}{T} \left(\sum_{i=1}^T h_i(x) \right) & \text{if Regression} \\ \text{majority_vote}(\{h_i(x)\}_{i=1}^T) & \text{if Classification} \end{cases}$$

有效性原因: 通过降低由于**不稳定学习器**(unstable learners)所造成的方差来降低错误率。不稳定学习器指那些在训练数据发生轻微变化的情况下, 学习结果会发生很大变化的学习器, 例如决策树。

Bagging 的优势:

- ①可以减小模型方差, 降低模型对训练数据的敏感性, 从而提高稳定性和泛化性, 常用于用来提高方差大但偏差小的基模型(决策树, 神经网络等单个模型不稳定的模型)的预测性能。
- ②便于**并行化**。多个抽样数据的获取及基模型的训练互相没有关联, 方便进行并行计算。

随机森林 Random Forest: “随机”是核心, “森林”是通过组合多棵决策树来构建模型。

假设使用三棵决策树组合成随机森林, 每棵树各不相同且预测结果相互独立, 每棵树的预测错误率为 40%。那么两棵树及以上预测错误的概率为: $0.4^3 + C_3^2 \times 0.4^2 \times (1 - 0.4) = 0.352$ 。

随机森林的算法原理:

随机森林在构建每棵树的时候, 为了**保证各棵树之间的独立性**, 通常会采用两到三层的随机性: 随机有放回的抽取样本, 随机选取 m 个特征, 随机选择特征取值进行分割(不遍历特征所有取值)。即:

使用 Bootstrap 抽样, 从训练集 D 中获得大小为 n 的抽样训练集 D_i 。从 D 个特征中随机选择 m 个特征, 基于 D_i 中随机选取的 m 个特征学习得到一个决策树 $h_i(x)$ 。

其中 m 的选取可以交叉验证, 也可以使用经验公式 $m = \log_2 D + 1$ 或 $m = \sqrt{D}$ 。

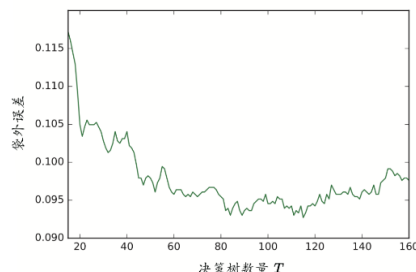
主要特点: 除了对样本进行有放回抽样, 还对特征进行随机抽样, 基本分类器最常见为决策树。

(**决策树局限性**: 局部最优, 分类边界「决策特征只涉及单个特征的逻辑判断, 导致决策边界是平行于坐标轴的直线, 这就限制了决策树对分类边界的表达能力」)

随机森林的性能评估:

分类间隔 margin: 正确分类某样本的决策树的比例减去错误分类样本决策树的比例, 此值越大越好。假设对样本 A 有 75% 的树分类正确, 那么分类间隔就是 $75\% - 25\% = 50\%$ 。

袋外误差 OOB (Out-Of-Bag Error): 随机森林对袋外样本(对某一棵树, 没有被抽样到训练集中的样本)的预测错误率。计算方式: 对每个样本, 计算把该样本作为袋外样本的树对该样本的分类情况, 以简单多数投票作为该样本的分类结果, 以误分样本个数占样本总数的比率作为随机森林的袋外误差。



随机森林与特征选择:

随机森林能够给出特征的重要性度量, 帮助进行特征选择。

计算特征的平均信息增益大小: 训练决策树时, 可以算出每个特征在每棵树中有多大的信息增益, 算出每个特征在随机森林中的平均信息增益, 作为该特征的重要性

计算每个特征对模型准确率的影响: 通过打乱样本中某一特征的特征值顺序, 产生新样本。新样本放入建立好的随机森林模型计算准确率, 对于重要的特征来说, 打乱顺序会极大的降低模型的准确率。

算法分析:

聚合模型并不改变单一模型的期望偏差, 但是可以缩小模型预测的方差。

设单棵决策树的方差为 $\sigma^2(x)$, 任两棵决策树之间的相关性为 $\rho(x)$, 随机森林一共由 T 棵决策树构成, 则模型方差 $\text{var}(x) = \rho(x)\sigma^2(x) + \frac{1-\rho(x)}{T}\sigma^2(x)$ 。证明见 hml 13.1。

随机森林优缺点:

优点:

能够处理很高维度的数据, 并且不用做特征选择。

对特征之间存在的多重共线性不敏感, 并且能够在一定程度上处理缺失数据和不均衡数据。在训练完后能够给出哪些特征比较重要。

容易做成并行化方法。

缺点:

处理噪音较大的小样本和低维数据集的问题上会过度拟合。

相对于决策树, 预测速度较慢, 模型可解释性较差。

随机森林其他功能——计算样本的相似度

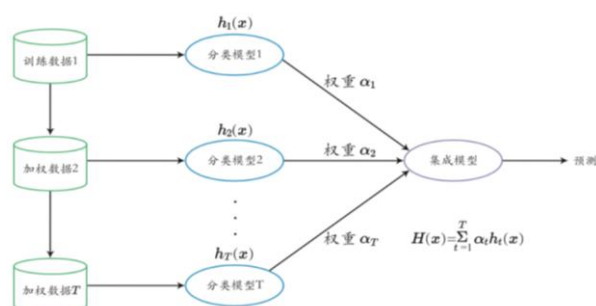
在建立随机森林的时候, 记录样本两两之间出现在同一叶子节点的次数, 生成相似性矩阵(proximity matrix)。如果越多的出现在同一叶节点, 说明这两个样本越相似; 如果越少的出现在同一叶节点, 说明两个样本差异越大。

Boosting: 起源于 PAC(Probably Approximately Correct)。串行训练多个模型, 后面的模型是基于前面模型的训练结果(误差)(重点关注前一个分错了的样本)。

$F_m(x) = F_{m-1}(x) + \arg \min \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + h(x_i))$ 典型

代表 AdaBoost。

在 AdaBoost 方法中, 将权重赋给每一个训练



样本, 迭代地学习 T 个分类器。初始权重均为 $\frac{1}{n}$, 每一步迭代, 被分错的样本的权重会降低。

第 t 个基础模型的**权重更迭方法**如下:

计算当前弱分类器误差(加权错误率, 弱分类器错误分类的样本的权重之和与训练集所有样本总权重之比): $\varepsilon_t = \frac{\sum_{i=1}^n \omega_i^t I(y_i \neq h_t(x_i))}{\sum_{i=1}^n \omega_i^t}$, 根据该误差计算权重: $\alpha_t = \frac{1}{2} \log(\frac{1-\varepsilon_t}{\varepsilon_t})$, 更新样

本权重: $\omega_i^{t+1} \leftarrow \frac{\omega_i^t e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$ (Z_t 为归一化因子), 最终的分类器为 $H(x) = \text{sign}(\sum_{t=1}^T \alpha_t h_t(x))$ 。

式中, I 为指示函数, 当 True 时值为 1, False 时值为 0。sign 函数在 $x \geq 0$ 时输出+1, 否则输出-1。当弱分类器对样本正确分类时 $y_i h_t(x_i) = 1$, 样本权重乘以小于1的因子 $e^{-\alpha_t}$, 相当于降低了正确分类的样本的权重。当弱分类器对样本错误分类时 $y_i h_t(x_i) = -1$, 样本权重乘以大于1的因子 $e^{-\alpha_t}$, 相当于提高了错误样本的权重。

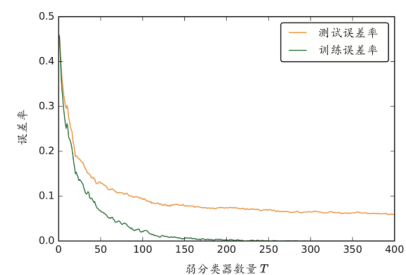
泛化误差对模型的测试集进行误差评估, 计算公式为

$\text{Error}_{\text{test}} = E_{x,y}(L(y_i, H(x_i)))$ 。理论上随着弱分类器数目 T 的增大, 泛化误差上界会增大。经验性结果表明, 即使 T 很大也不容易过度拟合。

AdaBoost 的优缺点:

优点: 防止过拟合, 减少方差, 减少偏差。

缺点: 由于集成了多个弱分类器, 模型的可解释性降低。当弱分类器太复杂或者效果太差时, 容易导致过拟合。对于异常值比较敏感。



神经网络

1. 神经元模型

McCulloch-Pitts 模型(M-P 模型):

对于输入信号 x_1, \dots, x_n 赋以权重 w_1, \dots, w_n , 计算其 $SUM = w_1x_1 + \dots + w_nx_n = W^T X$, 输出值 $Y = f(SUM - T)$ 。式中 W 是一个大小为 $s_l \times s_{l-1}$ 的矩阵, s_l, s_{l-1} 分别为本层、前一层的神经元数量。T 为阈值, f 为激活函数。将 $SUM - T$ 变为偏置项的形式 $WX + b$, 以 3 层神经网络为例, 输出应该是 $z = f(W^{(2)}f(W^{(1)}x + b^{(1)}) + b^{(2)})$, 也就是一个 2 层复合函数。

激活函数: Sigmoid: $f = \frac{1}{1 + e^{-x}} \in (0, 1)$, Tanh: $f = \tanh(x) \in (-1, 1)$, ReLU: $f = \max\{0, x\}$ 。

作用: 如果没有激活函数, 无论多少次复合, 神经网络的映射始终是线性函数, 无法处理非线性问题(仿射变换的叠加还是仿射变换)。**非线性变换相当于提升了数据的维度。**

二分类可选 **Logistic** 函数, 多分类可选 **softmax** 函数, 要求输出在 (a, b) 内的问题可选 $\frac{b-a}{2} \tanh(x) + \frac{b+a}{2}$ 。更一般的, 从 $[m, n]$ 均匀映射到 $[a, b]$ 的函数是 $f(x) = \frac{b-a}{n-m}(x-m) + a$,

均匀映射的含义是对于任意 $[u, v] \in [m, n]$, 都有 $\frac{v-u}{n-m} = \frac{f(v)-f(u)}{b-a}$ 。

Sigmoid 导数为 $f(x)(1-f(x)) \in (0, \frac{1}{4}]$, 很可能梯度消失, 有可能梯度爆炸。证明如下:

下: $y^{(i+1)} = f(z^{(i+1)}), z^{(i+1)} = w^{(i+1)}y^{(i)} + b^{(i)}$, BP 时 $\frac{\partial L}{\partial w^{(i)}} = \frac{\partial L}{\partial y^{(i)}} \frac{\partial y^{(i)}}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial w^{(i)}} = \frac{\partial L}{\partial y^{(i)}} f'(z^{(i)})y^{(i-1)}$,

式中 $\frac{\partial L}{\partial y^{(i)}} = \frac{\partial L}{\partial y^{(i+1)}} \frac{\partial y^{(i+1)}}{\partial z^{(i+1)}} \frac{\partial z^{(i+1)}}{\partial y^{(i)}} = \frac{\partial L}{\partial y^{(i+1)}} f'(z^{(i+1)})w^{(i+1)}$, 权重如果标准化, 则通常 $w^{(i+1)} \leq 1$,

此时有 $\left| \frac{\partial L}{\partial y^{(i)}} \right| \leq \frac{1}{4} \frac{\partial L}{\partial y^{(i+1)}}$ 。另外地, 当 $|w^{(i+1)}| > 4$ 时就有可能梯度爆炸。

ReLU 最简单, 但导致神经元在小于 0 时不会被更新。Leaky ReLU 给予负值 α 的斜率, RReLU 将 α 从均匀分布 $U(u, l)$ 中随机选取(测试时取均值 $\alpha = (u+l)/2$), PReLU 将 α 变为可训练参数。

Dice: $Dice(x) = p(x)x + (1-p(x))\alpha x$, $p(x) = \frac{1}{1 + \exp\left(-\frac{x - \mu_x}{\sqrt{\sigma_x^2 + \varepsilon}}\right)} = \sigma(BN(x))$, 式中 σ_x 代表

标准差, $\sigma(\cdot)$ 是 sigmoid 函数, 注意区分。Dice 使得拐点不再是固定的 0。

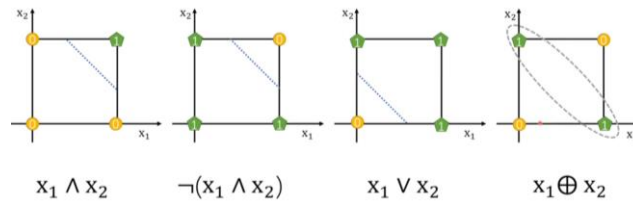
注: BN 是批归一化: $BN(x) = \gamma \frac{X - \mu}{\sqrt{\sigma^2 + \varepsilon}} + \beta$, 在 Dice 里 $\gamma=1, \beta=0$ 。BN 为了解决内

部协变量偏移 ICS 的问题, ICS 是指在神经网络训练过程中, 前一层参数的更新会导致后续层输入分布的改变, 每一层的输入分布不断变化, 需要不断适应新的输入分布, 这会减缓模型的收敛速度。式中 $\frac{X - \mu}{\sqrt{\sigma^2 + \varepsilon}}$ 是规范化, ε 是个小常数防止分母为零。 γ, β 用于缩放和平移规范化后的输出。

2. 感知机

感知机 Perceptron: 由两层神经元(输入层、输出层)组成。

能解决与、非问题，不能解决异或问题。why? 参考下图，依次为与、与非、或、异或(定义可参考 <https://blog.csdn.net/londa/article/details/119855241>):



如何通过感知机优化 M-P 模型中的 w, T ?

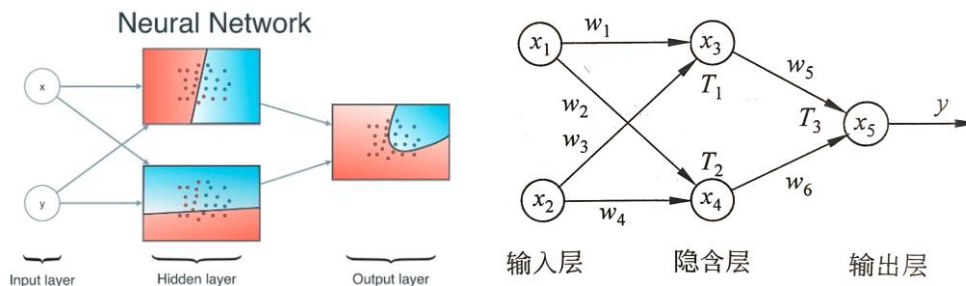
首先，将 T 视为 w_{n+1} 且 $w_{n+1} \equiv -1$ ，也就是将对阈值的学习转换为对权重的学习。接下来，对于训练样本 (x, y) ，设目前感知机输出的是 y ，权重调整为： $w_i \leftarrow w_i + \Delta w_i$ ，其中 $\Delta w_i = \eta(y - y)x_i$ ，式中学习率(Learning Rate) $\eta \in (0, 1)$ 。

3.(多层)神经网络 Neural Network

神经元模型只能处理线性问题，对于非线性问题的模拟，Neural Network 把多个神经元组合到一起，形成一个网络来运算，这便可以找出更复杂的数据规律。

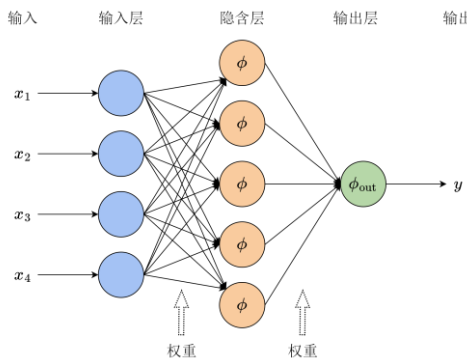
如左图，input layer 即输入层，hidden layer 即隐含层，output layer 即输出层。

如右图，我们依旧对每个 x_i 赋以权重 w_i ，对每个 SUM_i 赋以对应的 T_i ，得到对应的输出 $x_{3,4,5}$ ，其中 $x_{3,4}$ 继续用作输入， x_5 才是最终输出的 y 。



比如在 Iris 数据集中，输入层是 sepal length, sepal width, petal length, petal width，输出层是 setosa, versicolor, virginica。

当我们组合多个单层的感知机时，通常采用前馈 feedforward 结构，即将神经元分为不同的层，每一层只和其前后相邻层的神经元连接，层内以及间隔一层以上的神经元之间没有连接。前馈网络的层数是指权重的层数，将多个单层感知机按前馈结构组合起来，就形成了多层感知机 MLP (multi-layer perceptron)。下图就是一个两层的多层感知机：



神经网络体验网址: <https://playground.tensorflow.org/#activation=tanh&batchSize=10&dataset=circle®Dataset=reg-plane&learningRate=0.03®ularizationRate=0&noise=0&networkShape=5,6&seed=0.70553&showTestData=false&discretize=false&percTrainData=50&x=true&y=true&xTimesY=false&xSquared=false&ySquared=false&cosX=false&sinX=false&cosY=false&sinY=false&collectStats=false&problem=classification&initZero=false&hideText=false>

4.反向传播/误差逆传播算法

BP 算法(误差逆传播 Error Back Propagation)分为两个过程:

①正向传播:输入样本从输入层经隐含层至输出层,如果输出层的实际输出与样本标签(期望输出)不符,则转至误差反向传播,反之结束算法。

②反向传播:将输出误差(期望输出值减实际输出值)按原通路反向传递,经隐含层至输入层。反向传递过程中将误差分摊给各层的各个神经元,依据各神经元的误差通过梯度下降法调整神经元的 w, T 以减小误差。

反复执行①②直至达到预定训练次数或误差小于预设值。

不足:训练时间长、局部最优。

5.深度学习

深度学习即很深层的神经网络。

由于层数过多, BP 算法难收敛,因此有如下方法训练:

①无监督逐层训练:先预训练(按网络的层级结构对网络内部参数进行分组,对每组参数训练产生在局部较好的参数值),然后将这些局部最优的结果联合起来使用 BP 算法对整个网络参数微调。

②权值共享:让一组神经元使用相同的权重。比如 CNN 就使用卷积核滑动。

常见的深度学习方法:

卷积神经网络 CNN(Convolutional Neural Network):前向反馈神经网络,使用局部感受野,认为图像中局部范围内的像素之间联系较为紧密,距离较远的像素之间相关性较弱。所以每个神经元只对局部进行感知,然后在更高层将局部信息综合起来得到全局信息。

循环神经网络 RNN(Recurrent Neural Network):让上一层的输入作为下一层的输入。

生成式对抗网络 GAN(Generative Adversarial Network):生成模型 Generative Model 尽量生成真实的图片去欺骗判别模型 Discriminative Model,判别模型则尽量将生成图片和真实图片区分开来,形成一个动态的博弈过程。

自编码器 AE(AutoEncoder):通过编码器对输入信息降维(映射到低维空间,也称潜在空间),获取输入信息的表征(即表达方式)。解码器将降维的特征重新恢复到输入信息。编码器有 SAE、DAE、CAE、VAE。

深度置信网络 DBN(Deep Belief Nets):一种生成模型,通过训练权重 w ,让整个神经网络以最大概率生成训练数据。组成元件是受限玻尔兹曼机 RBM。RBM 是一种神经感知器,由一个显层(用于输入训练数据)和一个隐层(用作特征检测器)构成,二者神经元之间为双向全连接,内部无互连神经元。DBN 由 RBM 串联组成,上一个 RBM 的输出是下一个 RBM 的输入。

图卷积网络 GCN(Graph Convolution Network)GCN 是 CNN 在图数据上的推广,同时对节点特征信息和结构信息进行端对端学习。面对拓扑图(关系型数据)每个顶点的相邻顶点数都可能不同,没办法使用相同尺寸的卷积核进行卷积操作,需要使用空间域与谱域。空间域直接将卷积操作定义在每个节点的连接关系上,谱域是对图进行傅里叶变换后再卷积。

6.神经网络一些细节

①具有输入噪声的训练等价于 Tikhonov 正则化(L2)，由此得到**暂退法(dropout)**：在计算下一层之前将当前层中的一些节点置

零，置零概率为 p 。为了使 $E(h') = h$ ，

$$\text{让未置零的放大: } h' = \begin{cases} 0, & p \\ \frac{h}{1-p}, & \text{else} \end{cases}$$

当然，在测试集上一般不 dropout。

```
def dropout_layer(X, dropout):
```

```
    mask = (torch.rand(X.shape) > dropout).float()
```

```
    return mask * X / (1.0 - dropout)
```

通过上面代码可以将权重矩阵以 dropout 的概率暂退。

②**梯度爆炸** gradient exploding：参数更新过大，破坏了模型的稳定收敛。

梯度消失 gradient vanishing：参数更新过小，在每次更新时几乎不会移动，导致模型无法学习。

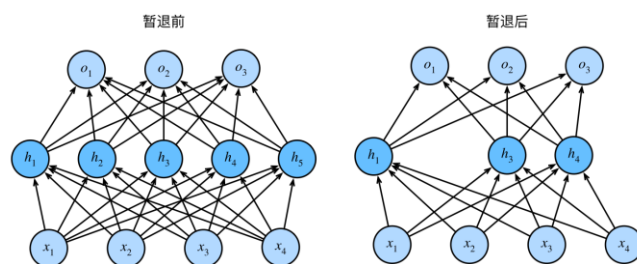
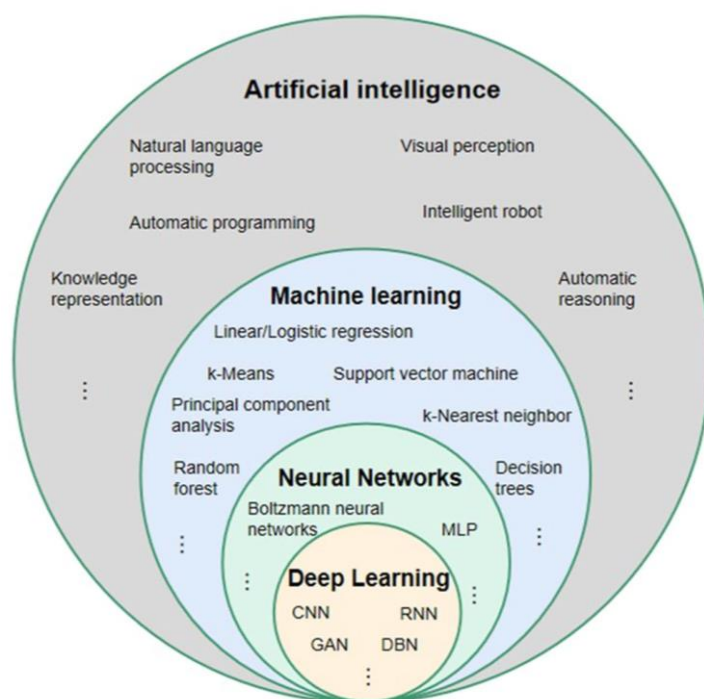


图4.6.1: dropout前后的多层感知机

7.机器学习&神经网络&深度学习的关系图



CNN

1.卷积核 convolution kernel/滤波器 filter

设输入大小 (n_h, n_w) ，卷积核大小 (k_h, k_w) ，则输出大小是 $(n_h - k_h + 1, n_w - k_w + 1)$ 。

如果填充 p_h 行 p_w 列，则输出大小是 $(n_h - k_h + p_h + 1, n_w - k_w + p_w + 1)$ 。

如果步幅是 s_h, s_w ，则输出大小是 $\left\lfloor \frac{n_h - k_h + p_h + s_h}{s_h} \right\rfloor, \left\lfloor \frac{n_w - k_w + p_w + s_w}{s_w} \right\rfloor$ 。

通常 $p_h = k_h - 1, p_w = k_w - 1$ 使得输入输出大小不变。如果 k_h 是偶数，可以在顶部填充 $\lceil p_h / 2 \rceil$ 行，在底部填充 $\lfloor p_h / 2 \rfloor$ 行。

感受野 receptive field：在前向传播期间可能影响 x 计算的所有元素(来自所有先前层)。

2.多输入多输出

多输入：对每个通道执行互相关操作，然后将结果相加。如右图， $(1 \times 1 + 2 \times 2 + 4 \times 3 + 5 \times 4) + (0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3) = 56$ 。

多输出：使用多个卷积核来生成多个输出通道，每个卷积核都会产生一个输出通道，因此一个卷积层可以生成多个不同的特征图。

也就是说，多输入是指对每个输入通道执行**单独**的卷积操作并将**结果相加**，多输出是指使用**多个**卷积核生成**多个输出通道**。

多输入多输出：对于每个输入通道执行多个卷积核的操作，从而产生多个输出通道。

$Y = X \cdot W + B$ ，其中各个元素的大小为：输入

$X: c_i \times n_h \times n_w$ ，核 $W: c_o \times c_i \times k_h \times k_w$ ，偏差

$B: c_o \times c_i$ ， $Y: c_o \times m_h \times m_w$ 。i:输入, o:输出, c:通道。

如右图，卷积核有 3 个输入通道和 2 个输出通道，分别对应输入、输出的 shape[0]。

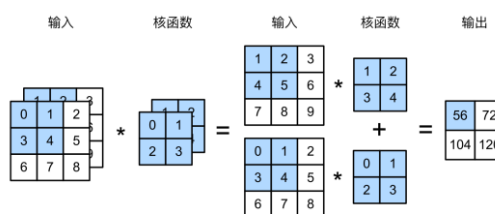
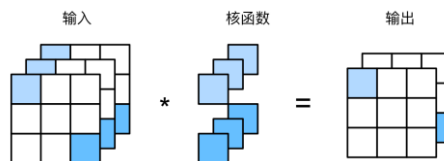


图6.4.1: 两个输入通道的互相关计算。



3.池化/汇聚 pooling

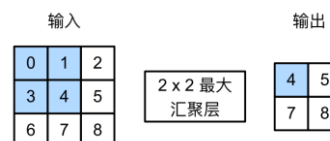
池化是一种降采样 downsampling 操作，用一个代表像素来替代一片区域内的所有像素，从而降低整体的复杂度。

汇聚层不包含参数，计算汇聚窗口中所有元素的最大值(maximum pooling)或平均值(average pooling)。

汇聚窗口形状为 $p \times q$ 的汇聚层被称为 $p \times q$ 汇聚层。

如右图， $\max(0, 1, 3, 4) = 4$ 。

汇聚层用于降低卷积层对位置的敏感性，同时降低对空间降采样表示的敏感性。



后续请参考：

<https://hml.boyuai.com/books/chapter9#93-%E7%94%A8%E5%8D%B7%E7%A7%AF%E7%A5%9E%E7%BB%8F%E7%BD%91%E7%BB%9C%E5%AE%8C%E6%88%90%E5%9B%BE%E5%83%8F%E5%88%86%E7%B1%BB%E4%BB%BB%E5%8A%A1>

RNN

之前都是默认数据来自于某种分布，并且所有样本都是独立同分布的。然而，例如，文章中的单词是按顺序写的，需要 RNN 处理这类序列信息。

PageRank 算法

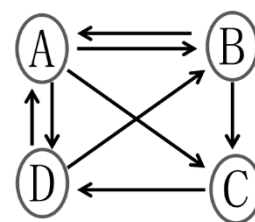
基本思想:

不只考虑词项，还考虑指向该网页的链接情况。如果一个网页被很多其他网页链接到，说明这个网页比较重要，PageRank 值较高。被越多优质(PageRank 高)的网页所指的网页，它是优质的概率就越大。

边上的权值表示网页之间随机跳转的概率(浏览者从当前网页点击到当前网页所指向网页的概率)。

如下例:

$$\text{转移矩阵 } M = \begin{pmatrix} 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & 0 & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} & 0 & 0 \\ \frac{1}{3} & 0 & 1 & 0 \end{pmatrix}, \text{初始 } T_0 \text{ 状态 Rank 值 } V_0 = \begin{pmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{pmatrix}。$$



第一轮: $T_0 \rightarrow T_1$

$$V_1 = MV = \begin{pmatrix} 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & 0 & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} & 0 & 0 \\ \frac{1}{3} & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \end{pmatrix} = \begin{pmatrix} \frac{1}{2} \cdot \frac{1}{4} + \frac{1}{2} \cdot \frac{1}{4} \\ \frac{1}{3} \cdot \frac{1}{4} + \frac{1}{2} \cdot \frac{1}{4} \\ \frac{1}{3} \cdot \frac{1}{4} + \frac{1}{2} \cdot \frac{1}{4} \\ \frac{1}{3} \cdot \frac{1}{4} + 1 \cdot \frac{1}{4} \end{pmatrix}$$

第二轮: $T_0 \rightarrow T_1 \rightarrow T_2$

$$V_2 = M^2V = \begin{pmatrix} 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & 0 & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} & 0 & 0 \\ \frac{1}{3} & 0 & 1 & 0 \end{pmatrix}^2 \begin{pmatrix} \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \end{pmatrix}$$

.....

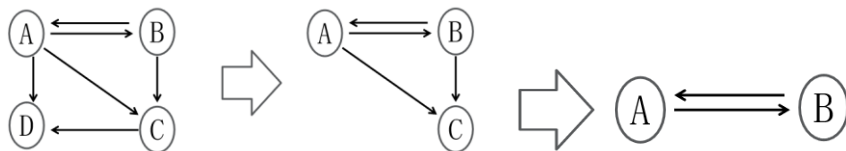
第 k 轮: $T_0 \rightarrow T_1 \rightarrow T_2 \rightarrow \cdots \rightarrow T_k$

$$V^k = M^k V$$

由于 M^k 是收敛的，所以迭代到一定次数后 V^k 趋于稳定。

PageRank 有效避免了词项作弊，但存在问题：

(1) **算法本身的问题：** Dead End 问题。这里可以把 D 节点去掉，但 C 也成了 dead end，所以 C 也要去掉。



(2) **作弊问题：**

① **链接农场问题：**

由互联网中的一部分网页组成，这些网页非常密集地互相连接在一起。通过创建一个堆砌大量链接而没有实质内容的网页，这些链接彼此互链或指向特定网站，以提高某个或者某些特定网页的 pagerank 值为目的。

② **交换链接：**

交换链接是指网站之间人为地互相增加对方网站的链接，是增加外链成本最低和使用最多的一种方法。由于交换链接可以增加网站的外链，提高网站的 pagerank 值。

③ **黄金链：**

指一些高权重的网站出售首页的链接给作弊网站，提高作弊网站的 pagerank 值。