

数据预处理

1. 在获取的数据中经常会遇到唯一属性，这些属性通常是人为添加的一些 `id` 属性，如存放在数据库中的自增的主键。

对于这些属性，它并不能刻画样本自身的分布规律。所以只需要简单的删除这些属性。

2. 对于数据中的某个属性，如果其方差很小，则意味着其识别能力较弱。极端情况下其方差为0，这意味着该属性在所有样本上的值都是恒定的。

因此可以设定一个阈值（如 10^{-3} ），将方差小于该阈值的属性直接剔除。

一、缺失值处理

1. 数据缺失值的产生的原因多种多样，主要分为客观原因和人为原因。

- 客观原因：比如数据存储的失败、存储器损坏、机械故障导致某段时间数据未能收集（对于定时数据采集而言）。
- 人为原因：由于人的主观失误（如：数据录入人员失误漏录了数据）、历史局限（如：数据在早期尚无记录）、有意隐瞒数据（如：在市场调查中被访人拒绝透露相关问题的答案，或者回答的问题是无效的）导致数据未能收集。

2. 缺失值的处理有三种方法：

- 直接使用含有缺失值的数据。

某些算法可以直接使用含有缺失值的情况，如决策树算法可以直接使用含有缺失值的数据。

- 优点：直接使用原始数据，排除了人工处理缺失值带来的信息损失。
- 缺点：只有少量的算法支持这种方式。

- 删除含有缺失值的数据。

最简单的办法就是删除含有缺失值的样本。

- 优点：简单、高效。
- 缺点：如果样本中的缺失值较少，则直接丢弃样本会损失大量的有效信息。这是对信息的极大浪费。

如果样本中包含大量的缺失值，只有少量的有效值，则该方法可以接受。

- 缺失值补全。

用最可能的值来插补缺失值。这也是在实际工程中应用最广泛的技术。

- 优点：保留了原始数据
- 缺点：计算复杂，而且当插补的值估计不准确时，会对后续模型引入额外的误差。

3. 缺失值补全常见有以下方法：

- 均值插补
- 同类均值插补
- 建模预测
- 高维映射
- 多重插补
- 压缩感知及矩阵补全

1.1 均值插补 && 同类均值插补

1. 均值插补：

- 如果样本的属性是连续值，则该属性的缺失值就以该属性有效值的平均值来插补。
- 如果样本的属性是离散值，则该属性的缺失值就以该属性有效值的众数（出现频率最高的值）来插补。

2. 均值插补在含有缺失值的属性上的所有缺失值都填补为同一个值。

而同类均值插补首先将样本进行分类，然后以该类中的样本的均值来插补缺失值。

1.2 建模预测

1. 建模预测的思想是：将缺失的属性作为预测目标，通过建立模型来预测。

2. 给定数据集 $\mathbb{D} = \{(\vec{x}_1, \tilde{y}_1), (\vec{x}_2, \tilde{y}_2), \dots, (\vec{x}_N, \tilde{y}_N)\}$ 。

假设属性 j 含有缺失值，根据 $x_{i,j}$ 是否缺失，将数据集划分为：

- $\mathbb{D}_1 = \{\vec{x}_i \mid x_{i,j} \neq null\}$ ：属性 j 有效的样本的集合。
- $\mathbb{D}_2 = \{\vec{x}_i \mid x_{i,j} = null\}$ ：属性 j 缺失的样本的集合。

将 \mathbb{D}_1 中的样本作为新的训练集，标签值重新定义为属性 j 的值，通过建模来完成属性 j 的学习。将 \mathbb{D}_2 中的样本作为测试集，通过学得模型来预测其属性 j 的值。

3. 这种方法的效果相对较好，但是该方法有个根本缺陷：

- 如果其他属性和属性 j 无关，则预测的结果无意义。
- 如果预测结果相当准确，则又说明属性 j 可以由其它属性计算得到，于是属性 j 信息冗余，没有必要纳入数据集中。

一般的情况是介于两者之间。

1.3 高维映射

1. 高维映射的思想是：将属性映射到高维空间。

2. 给定数据集 \mathbb{D} ，假设属性 j 的取值为离散值 $\{a_1, a_2, \dots, a_K\}$ 一共 K 个值，则将该属性扩展为 $K+1$ 个属性 $(j_1, j_2, \dots, j_{K+1})$ ，其中：

- 若样本在属性 j 上的取值为 a_k ，则样本在新的属性 j_k 上的取值为 1，在新的属性 $j_1, \dots, j_{k-1}, j_{k+1}, \dots, j_{K+1}$ 上的取值为 0。
- 若样本在属性 j 上缺失，则样本在新的属性 j_{K+1} 上的取值为 1，在新的属性 j_1, \dots, j_K 上取值为 0。

3. 对于连续特征，高维映射无法直接处理。可以在连续特征离散化之后，再进行高维映射。

4. 高维映射是最精确的做法，它完全保留了所有的信息，也未增加任何额外的信息。比如广告的 CTR 预估模型，预处理时会把所有变量都这样处理，达到几亿维。

- 优点：完整保留了原始数据的全部信息。
- 缺点：计算量大大提升。而且只有在样本量非常大的时候效果才好，否则会因为过于稀疏，效果很差。

1.4 多重插补

1. 多重插补 (Multiple Imputation: MI) 认为待插补的值是随机的，它的值来自于已观测到的值。

具体实践上通常是估计出待插补的值，然后再加上不同的噪声，形成多组可选插补值。然后根据某种选择依据，选取最合适的插补值。

2. 多重插补法的步骤：

- 通过变量之间的关系对缺失数据进行预测，利用蒙特卡洛方法生成多个完整的数据集。

- 在每个完整的数据集上进行训练，得到训练后的模型以及评价函数值。
- 对来自各个完整的数据集的结果，根据评价函数值进行选择，选择评价函数值最大的模型，其对应的插值就是最终的插补值。

1.5 压缩感知 && 矩阵补全

1. 在现实任务中，经常希望根据部分信息来恢复全部信息。压缩感知和矩阵补全就是用于完成这个任务。
2. 假定有长度为 n 的离散信号 \vec{x} 。根据奈奎斯特采样定理，当采样频率达到 \vec{x} 最高频率的两倍时，采样后的信号就保留了原信号的全部信息。

假定以远小于奈奎斯特采样定理要求的采样频率进行采样，得到了长度为 m 的采样后信号 \vec{y} ，其中 $m \ll n$ 。则有： $\vec{y} = \Phi \vec{x}$ 。

其中 $\Phi \in \mathbb{R}^{m \times n}$ 是对信号 \vec{x} 的测量矩阵，它确定了以什么样的频率采样以及如何将采样样本组成采样后的信号。

3. 通常在已知离散信号 \vec{x} 和测量矩阵 Φ 时要得到测量值 \vec{y} 很容易。但是如果给定测量值 \vec{y} 和测量矩阵 Φ ，要还原出原始信号 \vec{x} 比较困难。这是由于当 $m \ll n$ 时， $\vec{y} = \Phi \vec{x}$ 是一个欠定方程，无法简单的求出数值解。

假设存在某种线性变换 $\Psi \in \mathbb{R}^{n \times n}$ ，使得 $\vec{x} = \Psi \vec{s}$ ，其中 \vec{s} 也和 \vec{x} 一样是 n 维列向量，则有 $\vec{y} = \Phi \vec{x} = \Phi \Psi \vec{s}$ 。令 $A = \Phi \Psi \in \mathbb{R}^{m \times n}$ ，则 $\vec{y} = A \vec{s}$ 。

4. 如果能够从 \vec{y} 中恢复 \vec{s} ，则能够通过 $\vec{x} = \Psi \vec{s}$ 从 \vec{y} 中恢复出 \vec{x} 。

从数学意义上来看，这种做法没有解决任何问题。因为根据 $\vec{y} = A \vec{s}$ ，从 \vec{y} 中恢复 \vec{s} 这个问题仍然是欠定的。

但是在实际应用中发现，如果 \vec{s} 具有稀疏性（即大量的分量为零），则该问题能够很好地求解。这是因为稀疏性使得未知因素的影响大大减少。

此时 Ψ 称作稀疏基，而 A 的作用类似于字典，能够将信号转换为稀疏表示。

1.5.1 压缩感知

1. 与特征选择、稀疏表示不同，压缩感知侧重的是如何利用信号本身所具有的稀疏性，从部分观测样本中恢复原信号。
2. 压缩感知分为感知测量和重构恢复两个阶段。

- 感知测量：关注如何对原始信号进行处理以获得稀疏样本表示。常用的手段是傅里叶变换、小波变换、字典学习、稀疏编码等
- 重构恢复：关注的是如何基于稀疏性从少量观测中恢复原信号。

3. 限定等距性 Restricted Isometry Property: RIP：对于大小为 $m \times n, m \ll n$ 的矩阵 A ，若存在常数 $\delta_k \in (0, 1)$ ，使得对于任意向量 \vec{s} 和 A 的所有子矩阵 $A_k \in \mathbb{R}^{m \times k}$ ，都有：

$$(1 - \delta_k) \|\vec{s}\|_2^2 \leq \|A_k \vec{s}\|_2^2 \leq (1 + \delta_k) \|\vec{s}\|_2^2$$

则称 A 满足 k 限定等距性 k -RIP。

此时通过下面的最优化问题可以近乎完美的从 \vec{y} 中恢复出稀疏信号 \vec{s} ，进而恢复出 \vec{x} ：

$$\begin{aligned} \min_{\vec{s}} \quad & \|\vec{s}\|_0 \\ \text{s.t.} \quad & \vec{y} = A \vec{s} \end{aligned}$$

这里 L_0 范数表示向量中非零元素的个数。

4. 该最优化问题涉及 L_0 范数最小化，这是个 NP 难问题。但是 L_1 范数最小化在一定条件下与 L_0 范数最小化问题共解，于是实际上只需要求解最小化问题：

$$\begin{aligned} \min_{\vec{s}} \|\vec{s}\|_1 \\ s.t. \quad \vec{y} = \mathbf{A}\vec{s} \end{aligned}$$

可以将该问题转化为 LASSO 等价形式，然后通过近端梯度下降法来求解。

1.5.2 矩阵补全

1. 矩阵补全 matrix completion 解决的问题是：

$$\begin{aligned} \min_{\mathbf{X}} \text{rank}(\mathbf{X}) \\ s.t. \quad x_{i,j} = a_{i,j}, (i,j) \in \Omega \end{aligned}$$

其中

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,n} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \cdots & x_{m,n} \end{bmatrix} \quad \mathbf{A} = \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{bmatrix}$$

- \mathbf{A} 为观测矩阵，其中有很多缺失值。
- Ω 为 \mathbf{A} 中所有的有数值的下标的集合。
- \mathbf{X} 为需要恢复的稀疏信号， $\text{rank}(\mathbf{X})$ 为矩阵 \mathbf{X} 的秩。

该最优化问题也是一个 NP 难问题。

2. 考虑到 $\text{rank}(\mathbf{X})$ 在集合 $\{\mathbf{X} \in \mathbb{R}^{m \times n} : \|\mathbf{X}\|_F^2 \leq 1\}$ 上的凸包是 \mathbf{X} 的核范数 nuclear norm：

$$\|\mathbf{X}\|_* = \sum_{j=1}^{\min\{m,n\}} \sigma_j(\mathbf{X})$$

其中 $\sigma_j(\mathbf{X})$ 表示 \mathbf{X} 的奇异值。于是可以通过最小化矩阵核范数来近似求解：

$$\begin{aligned} \min_{\mathbf{X}} \|\mathbf{X}\|_* \\ s.t. \quad x_{i,j} = a_{i,j}, (i,j) \in \Omega \end{aligned}$$

该最优化问题是一个凸优化问题，可以通过半正定规划 Semi-Definite Programming:SDP 求解。

3. 理论研究表明：若 \mathbf{A} 的秩为 $r, n \leq m$ ，则只需要观察 $O(mr \log^2 m)$ 个元素就能够完美恢复出 \mathbf{A} 。

二、特征编码

2.1. 特征二元化

- 特征二元化的过程是将数值型的属性转换为布尔值的属性。通常用于假设属性取值为取值分布为伯努利分布的情形。
- 特征二元化的算法比较简单。对属性 j 指定一个阈值 ϵ 。
 - 如果样本在属性 j 上的值大于等于 ϵ ，则二元化之后为 1。
 - 如果样本在属性 j 上的值小于 ϵ ，则二元化之后为 0。
- 阈值 ϵ 是一个超参数，其选取需要结合模型和具体的任务来选择。

2.2. one-hot

1. 对于非数值属性，如 性别：[男，女]、国籍：[中国，美国，英国] 等等，可以构建一个到整数的映射。如 性别：[男，女] 属性中，将 男 映射为整数 1、女 映射为整数 0。

该方法的优点是简单。但是问题是，在这种处理方式中无序的属性被看成有序的。男 和 女 无法比较大小，但是 1 和 0 有大小。

解决的办法是采用独热码编码 One-Hot Encoding。

2. One-Hot Encoding 采用 N 位状态位来对 N 个可能的取值进行编码，每个取值都由独立的状态位来表示，并且在任意时刻只有其中的一位有效。

假设属性 j 的取值为非数值的离散集合 $\{a_1, a_2, \dots, a_K\}$ ，独热码将其扩展成 K 个属性，每个新属性代表属性 j 的一个状态位：若样本在属性 j 上的取值为 a_k ，则样本在新的属性 j_k 上的取值为 1，在新的属性 $j_1, \dots, j_{k-1}, j_{k+1}, \dots, j_K$ 上的取值为 0。

- 这种做法中，如果在 j_1, \dots, j_K 上取值全为 0，则表示发生了缺失。

缺失值也可以用第 $K + 1$ 个状态位来表示。

- 也可以扩展成 $K - 1$ 个属性，如果在 j_1, \dots, j_{K-1} 上取值全为 0，则表示样本在属性 j 上的取值为 a_K 。

3. One-Hot Encoding 的优点：

- 能够处理非数值属性。
- 在一定程度上也扩充了特征。如 性别 是一个属性，经过独热码编码之后变成了 是否男 和 是否女 两个属性。
- 编码后的属性是稀疏的，存在大量的零元分量。

4. 在决策树模型中，并不推荐对离散特征进行 one-hot。主要有两个原因：

- 产生样本切分不平衡的问题，此时且分增益会非常小。

如：国籍 这个离散特征经过独热码编码之后，会产生 是否中国、是否美国、是否英国、... 等一系列特征。在这一系列特征上，只有少量样本为 1，大量样本为 0。

这种划分的增益非常小，因为拆分之后：

- 较小的那个拆分样本集，它占总样本的比例太小。无论增益多大，乘以该比例之后几乎可以忽略。
- 较大的那个拆分样本集，它几乎就是原始的样本集，增益几乎为零。
- 影响决策树的学习。

决策树依赖的是数据的统计信息。而独热码编码会把数据切分到零散的小空间上。在这些零散的小空间上，统计信息是不准确的，学习效果变差。

本质是因为独热码编码之后的特征的表达能力较差的。该特征的预测能力被人为的拆分成多份，每一份与其他特征竞争最优划分点都失败。最终该特征得到的重要性会比实际值低。

2.3 离散化

1. 离散化用于将连续的数值属性转化为离散的数值属性。
2. 是否使用特征离散化，这背后是：使用“海量离散特征+简单模型”，还是“少量连续特征+复杂模型”。

- 对于线性模型，通常使用“海量离散特征+简单模型”。
 - 优点：模型简单。
 - 缺点：特征工程比较困难。但是一旦有成功的经验就可以推广，并且可以很多人并行研究。
- 对于非线性模型（如深度学习），通常使用“少量连续特征+复杂模型”。

- 优点是：不需要进行复杂的特征工程。
- 缺点是：模型复杂。

2.3.1 分桶

1. 离散化的常用方法是分桶。

- 将所有样本在连续的数值属性 j 的取值从小到大排列 $\{a_0, a_1, \dots, a_N\}$ 。
- 然后从小到大依次选择分桶边界 b_1, b_2, \dots, b_M 。其中：
 - M 为分桶的数量，它是一个超参数，需要人工指定。
 - 每个桶的大小 $b_{k+1} - b_k$ 也是一个超参数，需要人工指定。
- 给定属性 j 的取值 a_i ，对其进行分桶：
 - 如果 $a_i < b_1$ ，则分桶编号为 0。分桶后的属性的取值为 0。
 - 如果 $b_k \leq a_i < b_{k+1}$ ，则分桶编号为 k 。分桶后的属性的取值为 k 。
 - 如果 $a_i \geq b_M$ ，则分桶编号为 M 。分桶后的属性的取值为 M 。

2. 分桶的数量和边界通常需要人工指定。一般有两种方法：

- 根据业务领域的经验来指定。如：对年收入进行分桶时，根据2017年全国居民人均可支配收入约为 2.6 万元，可以选择桶的数量为5。其中：
 - 年收入小于 1.3 万元（人均的0.5倍），则为分桶 0。
 - 年收入在 1.3万元 ~ 5.2 万元（人均的0.5 ~ 2倍），则为分桶 1。
 - 年收入在 5.3万元 ~ 26万元（人均的2倍 ~ 10倍），则为分桶 2。
 - 年收入在 26万元 ~ 260万元（人均的10倍 ~ 100倍），则为分桶 3。
 - 年收入超过 260万，则为分桶 4。
- 根据模型指定。根据具体任务来训练分桶之后的数据集，通过超参数搜索来确定最优的分桶数量和分桶边界。

3. 选择分桶大小时，有一些经验指导：

- 分桶大小必须足够小，使得桶内的属性取值变化对样本标记的影响基本在一个不大的范围。
即不能出现这样的情况：单个分桶的内部，样本标记输出变化很大。
- 分桶大小必须足够大，使每个桶内都有足够的样本。
如果桶内样本太少，则随机性太大，不具有统计意义上的说服力。
- 每个桶内的样本尽量分布均匀。

2.3.2 特性

1. 在工业界很少直接将连续值作为逻辑回归模型的特征输入，而是将连续特征离散化为一系列 0/1 的离散特征。

其优势有：

- 离散化之后得到的稀疏向量，内积乘法运算速度更快，计算结果方便存储。
- 离散化之后的特征对于异常数据具有很强的鲁棒性。
如：销售额作为特征，当销售额在 `[30,100]` 之间时，为1，否则为 0。如果未离散化，则一个异常值 10000 会给模型造成很大的干扰。由于其数值较大，它对权重的学习影响较大。
- 逻辑回归属于广义线性模型，表达能力受限，只能描述线性关系。特征离散化之后，相当于引入了非线性，提升模型的表达能力，增强拟合能力。
假设某个连续特征 j ，它离散化为 M 个 0/1 特征 j_1, j_2, \dots, j_M 。则：

$w_j * x_j \rightarrow w_{j_1} * x'_{j_1} + w_{j_2} * x'_{j_2} + \dots + w_{j_M} * x'_{j_M}$ 。其中 $x'_{j_1}, \dots, x'_{j_M}$ 是离散化之后的新的特征，它们的取值空间都是 $\{0, 1\}$ 。

上式右侧是一个分段线性映射，其表达能力更强。

- 离散化之后可以进行特征交叉。假设有连续特征 j ，离散化为 N 个 0/1 特征；连续特征 k ，离散化为 M 个 0/1 特征，则分别进行离散化之后引入了 $M + N$ 个特征。

假设离散化时，并不是独立进行离散化，而是特征 j, k 联合进行离散化，则可以得到 $M \times N$ 个组合特征。这会进一步引入非线性，提高模型表达能力。

- 离散化之后，模型会更稳定。

如对销售额进行离散化，`[30, 100)` 作为一个区间。当销售额在40左右浮动时，并不会影响它离散化后的特征的值。

但是处于区间连接处的值要小心处理，另外如何划分区间也是需要仔细处理。

- 特征离散化简化了逻辑回归模型，同时降低模型过拟合的风险。

能够对抗过拟合的原因：经过特征离散化之后，模型不再拟合特征的具体值，而是拟合特征的某个概念。因此能够对抗数据的扰动，更具有鲁棒性。

另外它使得模型要拟合的值大幅度降低，也降低了模型的复杂度。

三、数据标准化、正则化

3.1. 数据标准化

- 数据标准化是将样本的属性取值缩放到某个指定的范围。

- 数据标准化的两个原因：

- 某些算法要求样本数据的属性取值具有零均值和单位方差。
- 样本不同属性具有不同量级时，消除数量级的影响。如下图所示为两个属性的目标函数的等高线。

- 数量级的差异将导致量级较大的属性占据主导地位。

从图中看到：如果样本的某个属性的量级特别巨大，将原本为椭圆的等高线压缩成直线，从而使得目标函数值仅依赖于该属性。

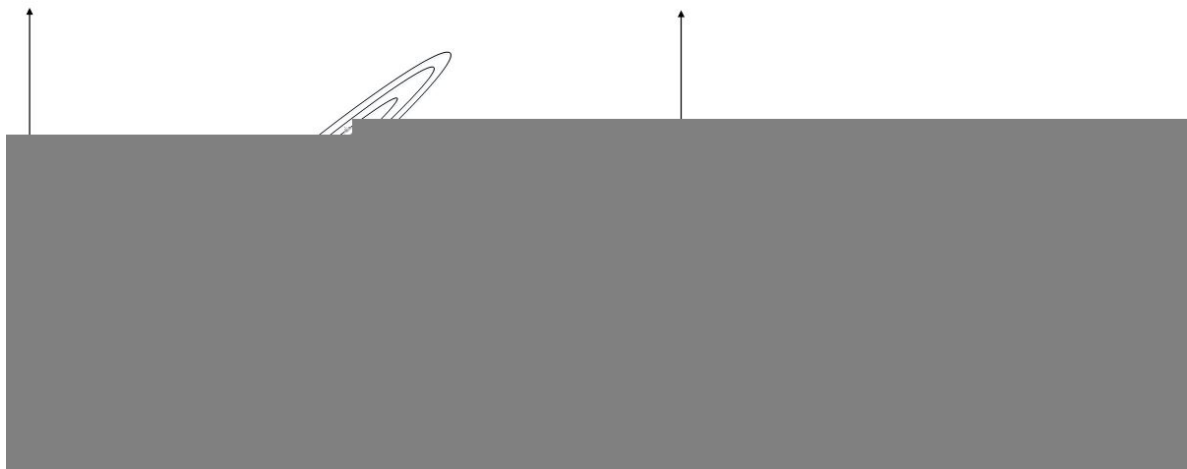
- 数量级的差异将导致迭代收敛速度减慢。

原始的特征进行梯度下降时，每一步梯度的方向会偏离最小值（等高线中心点）的方向，迭代次数较多，且学习率必须非常小，否则非常容易引起宽幅震荡。

标准化后进行梯度下降时，每一步梯度的方向都几乎指向最小值（等高线中心点）的方向，迭代次数较少。

- 所有依赖于样本距离的算法对于数据的数量级都非常敏感。

如 k 近邻算法需要计算距离当前样本最近的 k 个样本。当属性的量级不同时，选取的最近的 k 个样本也会不同。



3. 设数据集 $D = \{(\vec{x}_1, \tilde{y}_1), (\vec{x}_2, \tilde{y}_2), \dots, (\vec{x}_N, \tilde{y}_N)\}$, $\vec{x}_i = (x_{i,1}, \dots, x_{i,n})^T$ 。常用的标准化算法有：

- **min-max** 标准化：对于属性 j ，设所有样本在属性 j 上的最大值为 j_{\max} ，最小值为 j_{\min} 。则标准化后的属性值为：

$$\hat{x}_{i,j} = \frac{x_{i,j} - j_{\min}}{j_{\max} - j_{\min}}$$

标准化之后，所有样本在属性 j 上的取值都在 $[0, 1]$ 之间。

- **z-score** 标准化：对于属性 j ，设所有样本在属性 j 上的均值为 μ_j ，方差为 σ_j 。则标准化后的属性值为：

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sigma_j}$$

标准化之后，样本集的所有属性的均值都为 0，标准差均为 1。

4. 注意：如果数据集分为训练集、验证集和测试集，则：训练集、验证集、测试集使用相同标准化参数，该参数的值都是从训练集中得到。
- 如果使用 **min-max** 标准化，则属性 j 的标准化参数 j_{\max}, j_{\min} 都是从训练集中计算得到。
 - 如果使用 **z-score** 标准化，则属性 j 的标准化参数 μ_j, σ_j 都是从训练集中计算得到。

3.2. 数据正则化

1. 数据正则化是将样本的某个范数（如 L_1 范数）缩放到单位1。

设数据集 $D = \{(\vec{x}_1, \tilde{y}_1), (\vec{x}_2, \tilde{y}_2), \dots, (\vec{x}_N, \tilde{y}_N)\}$, $\vec{x}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,n})^T$ 。则样本 \vec{x}_i 正则化后的结果为：

$$\hat{\vec{x}}_i = \left(\frac{x_{i,1}}{L_p(\vec{x}_i)}, \frac{x_{i,2}}{L_p(\vec{x}_i)}, \dots, \frac{x_{i,n}}{L_p(\vec{x}_i)} \right)^T$$

其中 L_p 为范数： $L_p(\vec{x}_i) = (|x_{i,1}|^p + |x_{i,2}|^p + \dots + |x_{i,n}|^p)^{1/p}$ 。

2. 正则化的过程是针对单个样本的，对每个样本将它缩放到单位范数。

标准化是针对单个属性的，需要用到所有样本在该属性上的值。

3. 通常如果使用二次型（如点积）或者其他核方法计算两个样本之间的相似性时，该方法会很有用。

四、特征选择

1. 对于一个学习任务，给定了属性集，其中某些属性可能对于学习来说是很关键的，但是有些属性可能就意义不大。

- 对当前学习任务有用的属性称作相关特征 `relevant feature`。
- 对当前学习任务没有用的属性称作无关特征 `irrelevant feature`。

从给定的特征集合中选出相关特征子集的过程称作特征选择 `feature selection`。

2. 特征选择可能会降低模型的预测能力。因为被剔除的特征中可能包含了有效的信息，抛弃了这部分信息会一定程度上降低预测准确率。

这是计算复杂度和预测能力之间的折衷：

- 如果保留尽可能多的特征，则模型的预测能力会有所提升，但是计算复杂度会上升。
- 如果剔除尽可能多的特征，则模型的预测能力会有所下降，但是计算复杂度会下降。

4.1 特征选择原理

1. 特征选择是一个重要的数据预处理（`data preprocessing`）过程。在现实机器学习任务中，获取数据之后通常首先进行特征选择，然后再训练学习器。

进行特征选择的原因：

- 首先，在现实任务中经常会遇到维数灾难问题，这是由于属性过多造成的。如果能从中选择出重要的特征，使得后续学习过程仅仅需要在一部分特征上构建模型，则维数灾难问题会大大减轻。

从这个意义上讲，特征选择与降维技术有相似的动机。事实上它们是处理高维数据的两大主流技术。

- 其次，去除不相关特征往往会降低学习任务的难度。

2. 特征选择过程必须确保不丢失重要特征，否则后续学习过程会因为重要信息的缺失而无法获得很好的性能。

- 给定数据集，如果学习任务不同，则相关特征很可能不同，因此特征选择中的无关特征指的是与当前学习任务无关的特征。
- 有一类特征称作冗余特征 `redundant feature`，它们所包含的信息能从其他特征中推演出来。
 - 冗余特征在很多时候不起作用，去除它们能够减轻学习过程的负担。
 - 但如果冗余特征恰好对应了完成学习任务所需要的某个中间概念，则该冗余特征是有益的，能降低学习任务的难度。

这里暂且不讨论冗余特征，且假设初始的特征集合包含了所有的重要信息。

3. 要想从初始的特征集合中选取一个包含了所有重要信息的特征子集，如果没有任何领域知识作为先验假设，则只能遍历所有可能的特征组合。

这在计算上是不可行的，因为这样会遭遇组合爆炸，特征数量稍多就无法进行。

一个可选的方案是：

- 产生一个候选子集，评价出它的好坏。
- 基于评价结果产生下一个候选子集，再评价其好坏。
- 这个过程持续进行下去，直至无法找到更好的后续子集为止。

这里有两个问题：如何根据评价结果获取下一个候选特征子集？如何评价候选特征子集的好坏？

4.1.1 子集搜索

1. 如何根据评价结果获取下一个候选特征子集？这是一个子集搜索 `subset search` 问题。

2. 解决该问题的算法步骤如下：

- 给定特征集合 $\mathbb{A} = \{A_1, A_2, \dots, A_d\}$ ，首先将每个特征看作一个候选子集（即每个子集中只有一个元素），然后对这 d 个候选子集进行评价。
假设 A_2 最优，于是将 A_2 作为第一轮选定子集。
- 然后在上一轮的选定子集中加入一个特征，构成了包含两个特征的候选子集。
假定 A_2, A_5 最优，且优于 A_2 ，于是将 A_2, A_5 作为第二轮的选定子集。
-
- 假定在第 $k + 1$ 轮时，本轮的最优的特征子集不如上一轮的最优的特征子集，则停止生成候选子集，并将上一轮选定的特征子集作为特征选择的结果。

3. 这种逐渐增加相关特征的策略称作前向 forward 搜索。

类似地，如果从完整的特征集合开始，每次尝试去掉一个无关特征，这样逐渐减小特征的策略称作后向 backward 搜索。

4. 也可以将前向和后向搜索结合起来，每一轮逐渐增加选定相关特征（这些特征在后续迭代中确定不会被去除）、同时减少无关特征，这样的策略被称作双向 bidirectional 搜索。
5. 该策略是贪心的，因为它们仅仅考虑了使本轮选定集最优。但是除非进行穷举搜索，否则这样的问题无法避免。

4.1.2 子集评价

1. 如何评价候选特征子集的好坏？这是一个子集评价 subset evaluation 问题。
2. 给定数据集 \mathbb{D} ，假设所有属性均为离散型。对属性子集 \mathbb{A} ，假定根据其取值将 \mathbb{D} 分成了 V 个子集： $\{\mathbb{D}_1, \mathbb{D}_2, \dots, \mathbb{D}_V\}$
于是可以计算属性子集 \mathbb{A} 的信息增益：

$$g(\mathbb{D}, \mathbb{A}) = H(\mathbb{D}) - H(\mathbb{D} | \mathbb{A}) = H(\mathbb{D}) - \sum_{v=1}^V \frac{|\mathbb{D}_v|}{|\mathbb{D}|} H(\mathbb{D}_v)$$

其中 $|\cdot|$ 为集合大小， $H(\cdot)$ 为熵。

信息增益越大，则表明特征子集 \mathbb{A} 包含的有助于分类的信息越多。于是对于每个候选特征子集，可以基于训练数据集 \mathbb{D} 来计算其信息增益作为评价准则。

3. 更一般地，特征子集 \mathbb{A} 实际上确定了对数据集 \mathbb{D} 的一个划分规则。
- 每个划分区域对应着 \mathbb{A} 上的一个取值，而样本标记信息 y 则对应着 \mathbb{D} 的真实划分。
 - 通过估算这两种划分之间的差异，就能对 \mathbb{A} 进行评价：与 y 对应的划分的差异越小，则说明 \mathbb{A} 越好。
 - 信息熵仅仅是判断这个差异的一种方法，其他能判断这两个划分差异的机制都能够用于特征子集的评价。
4. 将特征子集搜索机制与子集评价机制结合就能得到特征选择方法。
- 事实上，决策树可以用于特征选择，所有树结点的划分属性所组成的集合就是选择出来的特征子集。
 - 其他特征选择方法本质上都是显式或者隐式地结合了某些子集搜索机制和子集评价机制。
5. 常见的特征选择方法大致可分为三类：过滤式 filter、包裹式 wrapper、嵌入式 embedding。

4.2 过滤式选择

1. 过滤式方法先对数据集进行特征选择，然后再训练学习器，特征选择过程与后续学习器无关。
这相当于先用特征选择过程对初始特征进行过滤，再用过滤后的特征来训练模型。

2. **Relief: Relevant Features** 是一种著名的过滤式特征选择方法，该方法设计了一个相关统计量来度量特征的重要性。

- 该统计量是一个向量，其中每个分量都对应于一个初始特征。特征子集的重要性则是由该子集中每个特征所对应的相关统计量分量之和来决定的。
- 最终只需要指定一个阈值 τ ，然后选择比 τ 大的相关统计量分量所对应的特征即可。

也可以指定特征个数 k ，然后选择相关统计量分量最大的 k 个特征。

3. 给定训练集 $\mathbb{D} = \{(\vec{x}_1, \tilde{y}_1), (\vec{x}_2, \tilde{y}_2), \dots, (\vec{x}_N, \tilde{y}_N)\}$, $\tilde{y}_i \in \{0, 1\}$ 。对于每个样本 \vec{x}_i ：

- **Relief** 先在 \vec{x}_i 同类样本中寻找其最近邻 \vec{x}_{nh_i} ，称作猜中近邻 **near-hit**。
- 然后从 \vec{x}_i 的异类样本中寻找其最近邻 \vec{x}_{nm_i} ，称作猜错近邻 **near-miss**。
- 然后相关统计量对应于属性 j 的分量为：

$$\delta_j = \sum_{i=1}^N \left(-\text{diff}(x_{i,j}, x_{nh_i,j})^2 + \text{diff}(x_{i,j}, x_{nm_i,j})^2 \right)$$

其中 $\text{diff}(x_{a,j}, x_{b,j})$ 为两个样本在属性 j 上的差异值，其结果取决于该属性是离散的还是连续的：

- 如果属性 j 是离散的，则：

$$\text{diff}(x_{a,j}, x_{b,j}) = \begin{cases} 0, & \text{if } x_{a,j} = x_{b,j} \\ 1, & \text{else} \end{cases}$$

- 如果属性 j 是连续的，则：

$$\text{diff}(x_{a,j}, x_{b,j}) = |x_{a,j} - x_{b,j}|$$

注意：此时 $x_{a,j}, x_{b,j}$ 需要标准化到 $[0, 1]$ 区间。

4. 从公式