



数据预处理：数据规约-维度规约

34

□ 基本计算思路 ---- PCA

- 1.对每个样本(N 维, M 个样本)提取出属性组成一个数字向量
- 2.对所有样本里的每个属性的取值进行归一化（按属性归一化），以消除不同属性的取值范围等不同带来的影响，得到一个 $N \times M$ 样本矩阵 X (归一化)；
- 3.该矩阵 X 乘以该矩阵的转置为协方差矩阵，这个协方差矩阵是可对角化的，对角化后剩下的元素为特征值，每个特征值对应一个特征向量；
- 4.选取最大的 K 个特征值（其中 K 即为PCA的主元（PC）数， K 越少，越降低数据量，但信息丢失也越大，识别效果也越差），将这 K 个特征值对应的特征向量（特征向量要标准化）组成新的矩阵($N \times K$)；
- 5.将新的矩阵转置($K \times N$)后乘以样本向量($N \times M$)即可得到降维后的数据（这些数据是原数据中相对较为主要的，而数据量 $K \times M$ 一般也远远小于原数据量 $N \times M$ ）。



数据预处理：数据规约-维度规约

35

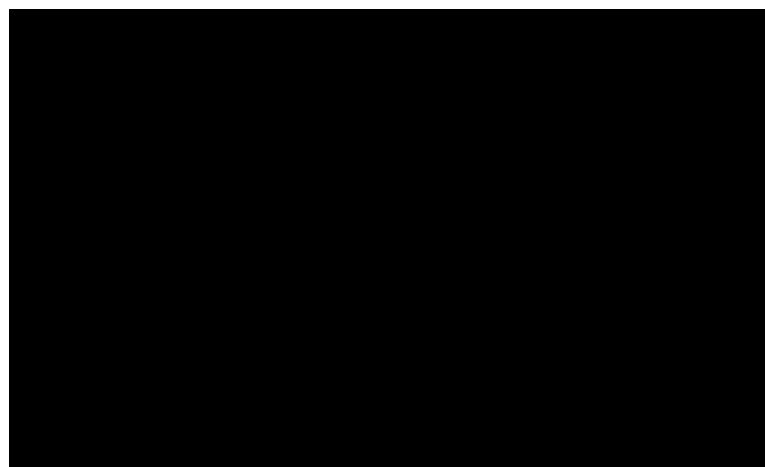
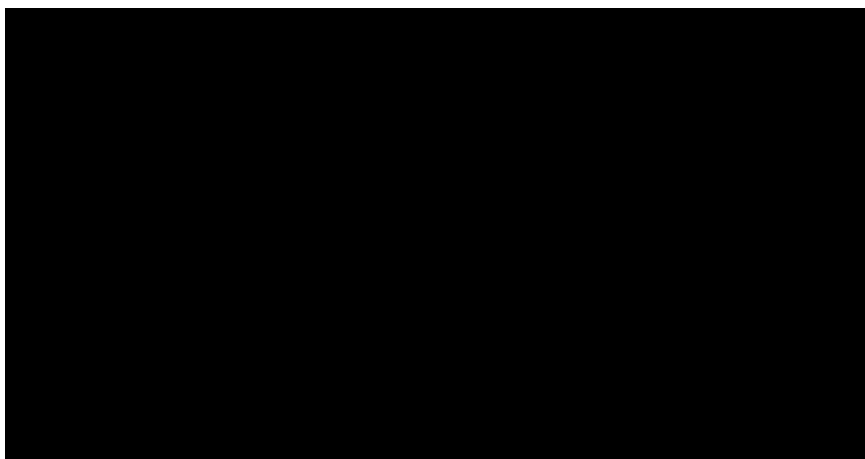
样本矩阵(10个城市样本, 8个属性)的转置 X^T



数据归一化并得
到协方差矩阵



三个主成分 (8*3)





数据预处理：数据规约-维度规约

36

先把大图像分成 16×16 （256）的小图像块，把小图像块当成一个256维的向量，所有256维向量拼接成新的数据矩阵，对其进行归一化和PCA压缩（取前四个特征值，取前八个，前16个，一直到前256个特征值），压缩完以后需要重构图像就会得到该效果





数据预处理：数据规约-数值规约

37

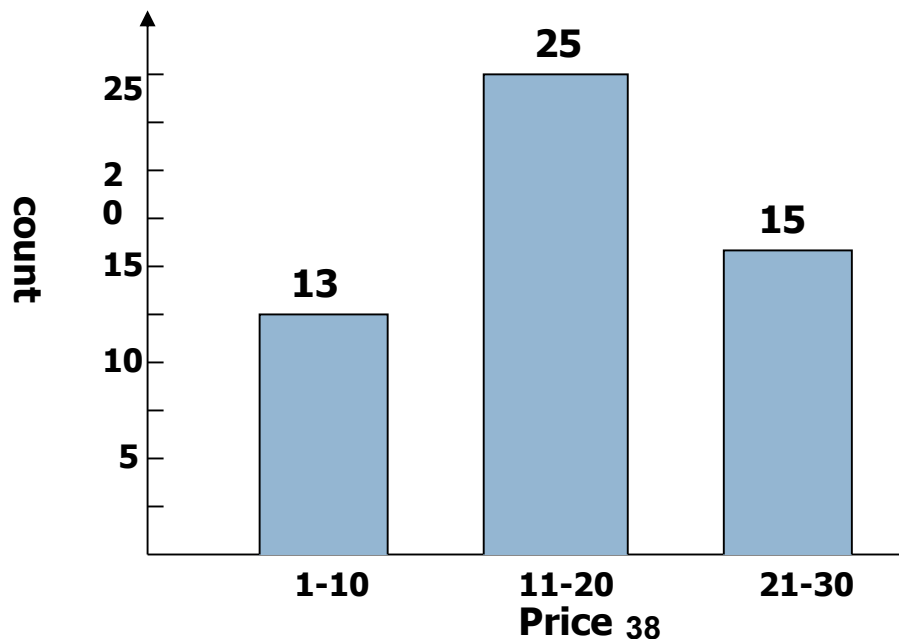
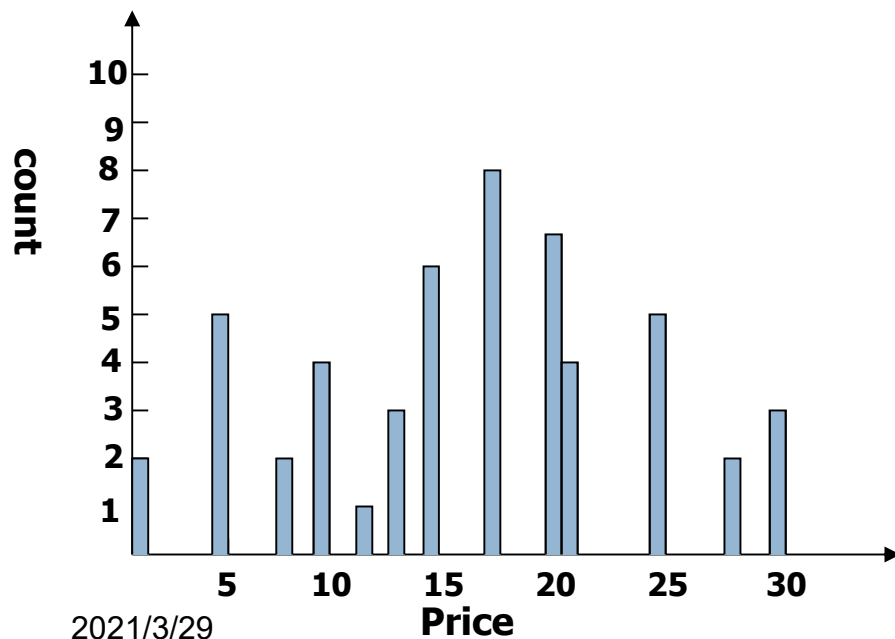
- 通过选择替代的、较小的数据表示形式来减少数据量
- 有参方法
 - 使用一个参数模型估计数据，最后只要存储参数即可，不用存储数据（除了可能的离群点）
 - 常用方法：线性回归方法；多元回归；对数线性模型；
- 无参方法
 - 不使用模型的方法存储数据
 - 常用方法：直方图，聚类，抽样



直方图

38

- 类似于分箱技术，是一种流行的数据归约方式
- 将属性值划分为不相交的子集，或“桶”
- 桶安放在水平轴上，而桶的高度（和面积）是该桶所代表的值的平均频率。
- 每个桶只表示单个属性值，则称其为“单桶”。通常，“桶”表示给定属性的一个连续空间





资料推荐

39

- 数据挖掘导论 第2章：数据，人民邮电出版社
- 数据挖掘原理与算法 第2章， 清华大学出版社
- T.C. Redman *Data Quality: The Field Guide*. January 2001
- I.T.Jolliffe. *Principal Component Analysis*. Springer Verlag, 2nd edition, October 2002.
- *Feature selection algorithms: A survey and experimental evaluation*, ICDM 2003



特征工程

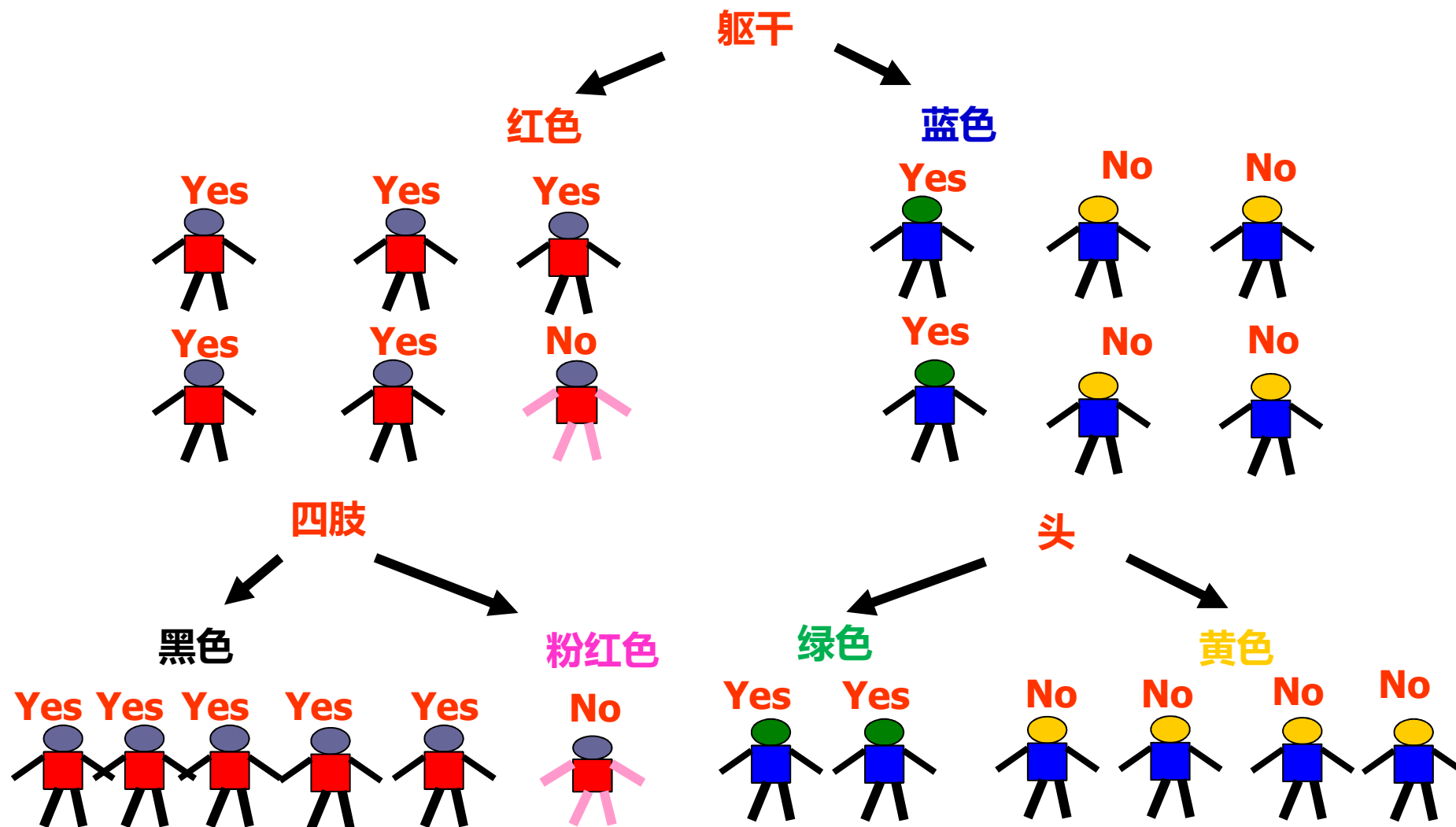
40

- 特征工程定义
- 特征工程的流程
- 特征学习
- 案例学习
- 参考文献



特征工程

41



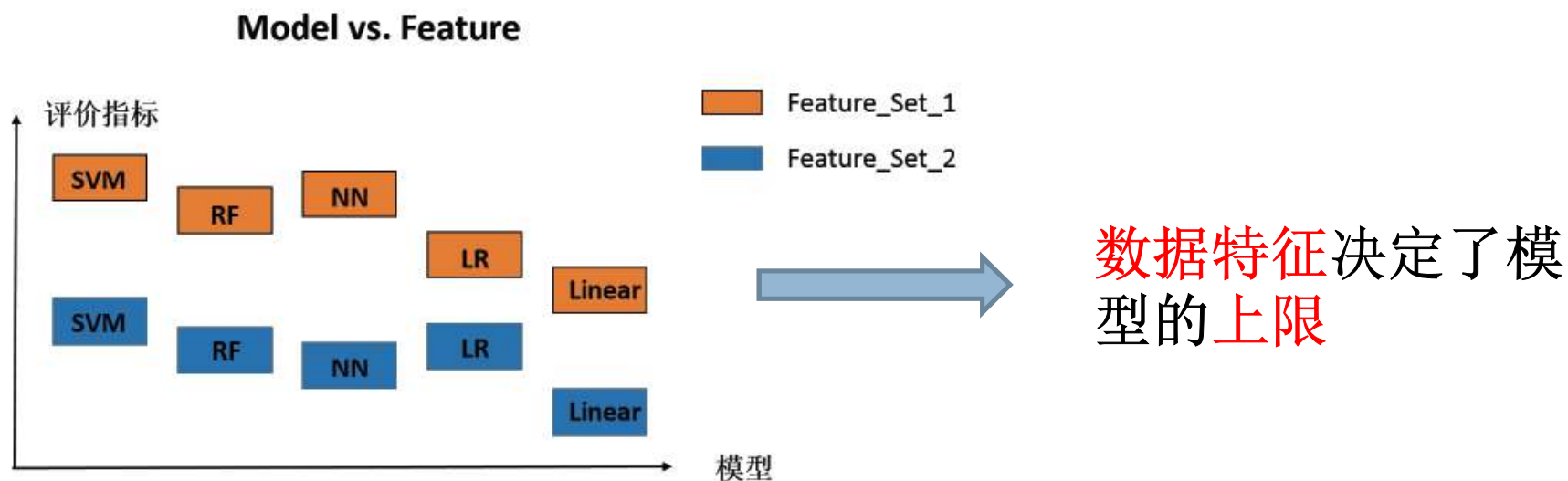


特征工程

42

□ 特征工程是什么？

在数据预处理以后（或者数据预处理过程中），如何从数据中提取有效的特征，使这些特征能够尽可能的表达原始数据中的信息，使得后续建立的数据模型能达到更好的效果，就是特征工程(Feature Engineering)所要做的工作。



- Feature 决定模型 UpperBound
- Model 决定接近 UpperBound的程度
- 不同的问题下Model的表现的不同



特征工程

43

□ 特征工程的意义

著名数据科学家Andrew Ng 对特征工程这样描述的：“虽然提取数据特征是非常困难、耗时并且需要相关领域的专家知识，但是机器学习应用的**基础**就是特征工程”

□ 特征越好，灵活性越强

好的特征能使一般的模型也能获得很好的性能，在不复杂的模型上运行速度很快，并且容易理解和维护。

□ 特征越好，构建的模型越简单

好的特征不需要花太多的时间去寻找最优参数，降低了模型的复杂度，使模型趋于简单。

□ 特征越好，模型的性能越出色

好的特征能够使模型表现越出色是毫无疑问的，提升模型的性能。

如何去做特征工程？

目的就是提



特征工程的流程

44

□ 特征工程（重复迭代）的流程

1. 对特征进行头脑风暴

深入分析问题，观察数据的基本统计信息，结合问题的相关领域知识和参考其他问题的相关特征工程的方法并应用到自身的问题中来。

2. 特征的设计

人工设计特征、自动提取特征，或者将两者相互结合，得到模型中所使用的特征。

3. 特征的选择

使用不同的特征重要性评分方法或者特征选择方法，对特征的有效性进行分析，选出有效的特征。

4. 评估模型

利用所选择的特征对测试数据进行预测，评估模型的性能。

5. 上线测试

通过在线测试的效果判断特征是否有效，若不能达到要求，则重复2-5步骤，直到模型的性能达到要求。



特征的设计

45

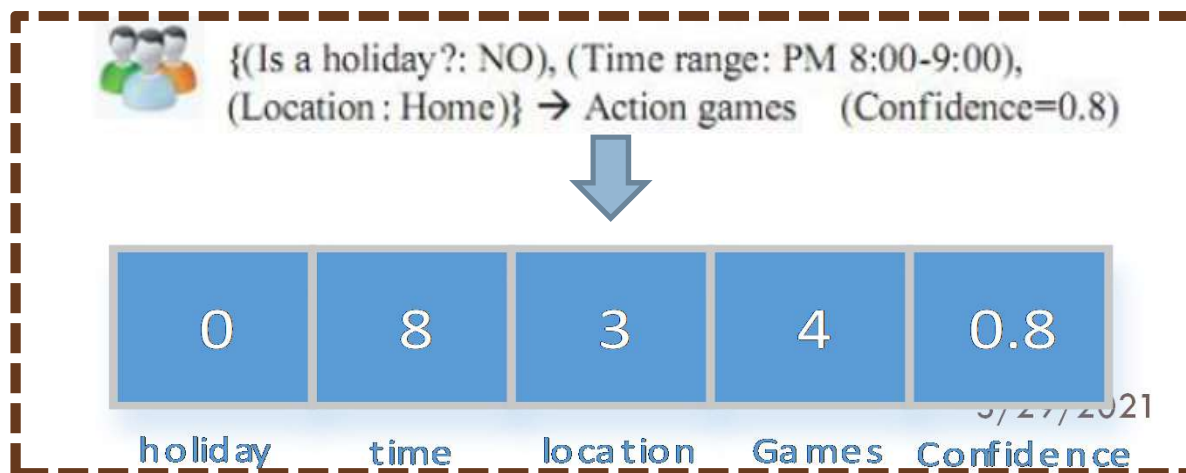
□ 从原始数据中如何设计特征？

□ 基本特征的提取

基本特征的提取过程就是对原始数据进行**预处理**将其转化成可以使用的数值特征。常见的方法有：数据的归一化、离散化、缺失值补全和数据变化等。

□ 创建新的特征

根据对应的领域知识，在基本特征的基础上进行特征之间的**比值**和**交叉变化**来构建新的特征。





特征的设计

46

- 从原始数据中如何设计特征？
 - 独热特征表示 One-hot Representation
 - 将每个属性表示成一个很长的向量（每维代表一个属性值，如词语）
 - 函数： $[0, 0, 1, 0, 0, \dots, 0, 0, 0, 0]$
 - 图像： $[0, 0, 0, 0, 0, \dots, 0, 0, 0, 1]$
 - 优势：简洁明了，缺陷：
 - “维度灾难” 问题：尤其是我们所构建的语料库包含的词语数目非常多的时候，独热表征在空间和时间上的开销都是十分巨大的
 - “语义鸿沟” 现象：任意两个词之间都是完全孤立的，是无法刻画句子中词语的语序信息的（之前提到的词袋模型也是如此）。例如，我们是无法通过独热表征来判断“函数”与“偶函数”之间的联系（但实际上这两个词语是非常相关的）。

3/29/2021

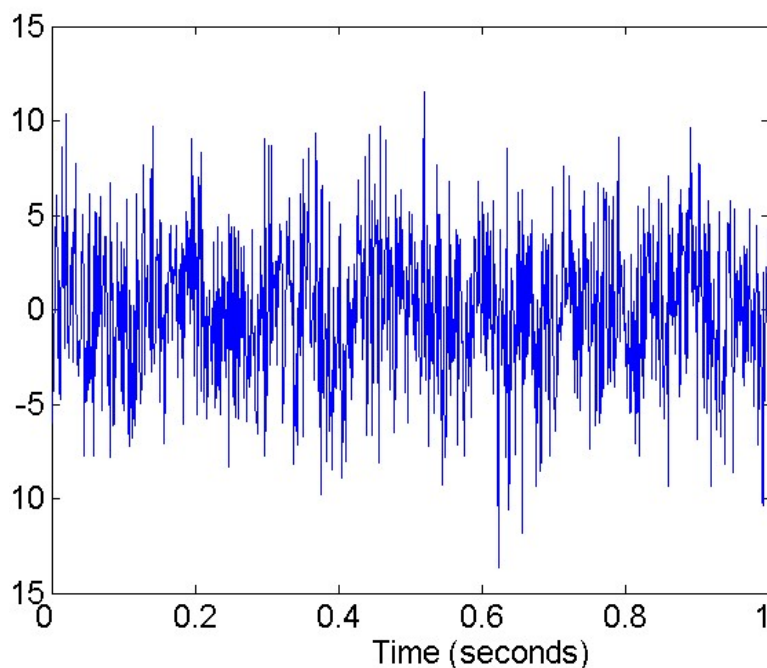


特征的设计

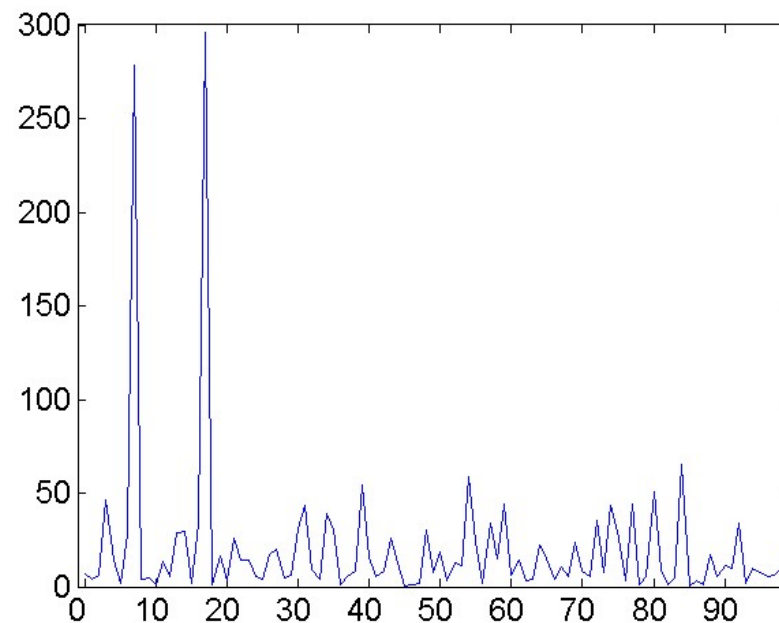
•47

□ 从原始数据中如何设计特征？

- 左图是根据两个Sin函数（分别是每秒7个和17个周期），以及一些噪声数据得到的序列图，右图是由傅立叶变换得到了频率图，可以看出变换后成功得到了两个概率最大的频率7和17（其中纵坐标是振幅，即概率值）



Two Sine Waves(正弦波) + Noise



Frequency



特征的设计

•48

埃塞俄比亚总理获得诺贝尔和平奖 特朗普愿望又一次落空

2019年10月11日消息，挪威诺贝尔委员会今天宣布，2019诺贝尔和平奖授予埃塞俄比亚总理阿比·艾哈迈德·阿里，表彰他“为实现和平与国际合作所作的努力，以及他为解决与邻国厄立特里亚的边界冲突所采取的果断举措”。

□ 从原始数据中如何设计特征？

□ TF-IDF（词频-逆文档率）

□ 算法简单高效，常被工业界用于最开始的数据预处理步骤

□ 主要思想：找到能代表该句子中的“关键词”（以文本数据为例）

□ 词频（TF, Term Frequency）

■ **TF = 某个词（特征值）在句子（数据）中出现的次数**

□ 逆文档率（IDF, Inverse Document Frequency）

■ **IDF = $\log(\text{语料库（数据库）的句子（数据）总数} / \text{包含该词（特征值）的句子（数据）总数})$**

□ 每个特征值（词）的重要性

□ **$w_{ij} = \text{tf} * \text{idf} = \text{TF}_{ij} * \log(N / \text{DF}_i)$**

可能会有
微小变形



特征的设计

•49

埃塞俄比亚总理获得诺贝尔和平奖 特朗普愿望又一次落空

2019年10月11日消息，挪威诺贝尔委员会今天宣布，2019诺贝尔和平奖授予埃塞俄比亚总理阿比·艾哈迈德·阿里，表彰他“为实现和平与国际合作所作的努力，以及他为解决与邻国厄立特里亚的边界冲突所采取的果断举措”。

□ 从原始数据中如何设计特征？

□ TF-IDF（词频-逆文档率）

□ 每个特征值（词）的重要性

□ $w_{ij} = tf * idf = TF_{ij} * \log(N/DF_i)$

■ 如何找到关键特征（词）？

- ① 根据 **TF** 可以找到并删去一个句子中的高频词（特征值）（比如停用词，“的”，“是”，“了”等）
- ② 根据 **IDF** 继续对句子中剩下的词进行权重赋值并排序，在数据库中越常见的词（特征值）权重越小
- ③ 根据 **TF-IDF** 我们可以得到一个句子（数据）中所有词（特征值）的 **TF-IDF** 值，进而排序筛选得到每个句子最有代表性的特征（“关键词”）



特征的设计

•50

□ 从原始数据中如何设计特征？

□ TF-IDF (词频-逆文档率) $w_{ij} = \text{tf} * \text{idf} = \text{TF}_{ij} * \log(N/\text{DF}_i)$

□ d_1 (A, B, C, C, S, D, A, B, T, S, S, S, T, W, W)

□ d_2 (C, S, S, T, W, W, A, B, S, B)

文档中关键词
总数=25

正则化的TF

	d_1	d_2
A	0.08	0.04
B	0.08	0.08
C	0.08	0.04
D	0.04	0.00
S	0.16	0.12
T	0.08	0.04
W	0.08	0.08

IDF

	$\ln(((1+ D)/ D_t))$
A	0.4
B	0.4
C	0.4
D	1.1
S	0.4
T	0.4
W	0.4

TF-IDF

	d_1	d_2
A	0.032	0.016
B	0.032	0.032
C	0.032	0.016
D	0.044	0.000
S	0.064	0.048
T	0.032	0.016
W	0.032	0.032



特征的设计

•51

- 从原始数据中如何设计特征？
 - TF-IDF（词频-逆文档率） $w_{ij} = \text{tf} * \text{idf} = \text{TF}_{ij} * \log(N/\text{DF}_i)$
- 优点
 - 简单快速的词（特征）重要性表示方法，结果比较符合实际情况
 - 应用广泛：不仅限于文本数据
- 缺点
 - 单纯以“词频”衡量一个词的重要性，不够全面，有时重要的词可能出现次数并不多
 - 无法体现词的位置信息、顺序信息，出现位置靠前的词与出现位置靠后的词，都被视为重要性相同
 - 无法发现词（特征）的隐含联系，如同义词等



特征的设计

52

- 举例：第二届“中国高校计算机大赛-大数据挑战赛”
- 赛题描述/数据：<http://bdc.saikr.com/vse/bdc/2017>
- 简单的说，该赛题的求解目标是利用数据分析将人工的鼠标轨迹和代码生成的鼠标轨迹区分开来。这里的鼠标轨迹是指一种完成一种验证手段——拖动滑块到指定区域时鼠标的轨迹。



- 原始数据格式：一系列连续点的坐标及其对应时间，目标点的坐标

例如(2,3,4),(2,5,6)(4,3,7) (4,3)，该轨迹中含有三个点的坐标，以(x,y, time)的时间表示，终点坐标为(4,3)

3/29/2021



特征的设计

53

□ 从原始数据中如何设计特征？

□ 基本特征的提取

- 轨迹运动数据的统计值，如运动速度/加速度/角加速度/角速度的均值/极值/最值/中位数 等等
- 轨迹的描述：运动在x轴方向是否为单向，曲线平滑程度， 等等

□ 创建新的特征

- 基本特征的简单二元运算， 加/减/乘/除/平方和/和平方/倒数和
- 运动数据在某一维上的偏导
- 领域专家知识

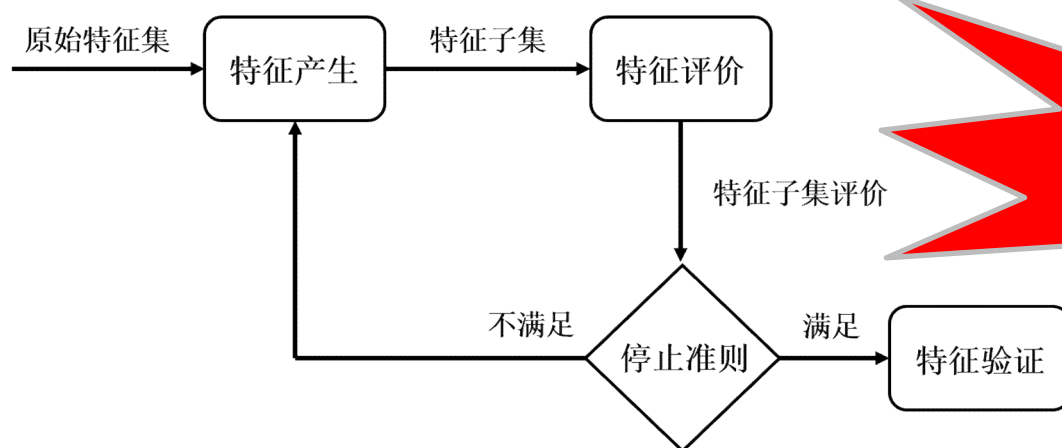


特征的选择

54

□ 如何挑选有效的特征（**Subset Selection问题**）？

在实际应用中，特征的数量往往比较多，其中可能会存在不相关的特征。而特征数量越多，分析特征、训练模型所需要的时间就越长，同时容易引起“维度灾难”，使得模型更加复杂。**特征选择**通过剔除不相关的特征或冗余的特征来**减少特征数量**，从而简化了模型并且提升了模型的泛化能力。



从特征中挑选
有效的
特征子集

特征选择的过程

3/29/2021



特征子集产生过程

55

□ □ 如何生成特征子集？

特征选择的本质上是一个组合优化的问题，求解组合优化问题最直接的方法就是搜索.根据不同的搜索策略，可以将搜索的算法分为完全搜索(Complete), 启发式搜索(Heuristic) 和随机搜索(Random) 三大类。

1. 采用全局最优搜索策略的过程产生方法

全局最优搜索策略可以分为穷举搜索与非穷举搜索两类。穷举搜索策略有遍历所有特征和以广度优先搜索的策略，这两种搜索策略都枚举了所有的特征组合，复杂度为 2^n 。

2. 采用启发式搜索策略的过程产生方法

启发式搜索的基本思想是增加关于要解决问题的解某些特征，以便指导搜索向最有希望的方向发展。启发式搜索是搜索是在搜索的最优性和计算量之间做一个折中的搜索策略。

3. 采用随机算法搜索策略的产生方法

特征选择本质上是一个组合优化问题，求解这类问题可采用非全局最优目标搜索方法，其实现是依靠带一定智能的随机搜索策略。(如模拟退火，遗传算法等)



特征子集产生过程

56

□ 举例：

- 对于前面提到的比赛的初步数据特征，我们采用了step-forward的方法来粗糙地选出200个特征。
- step-forward本质上是一种基于贪心策略的搜索方法，每次从未被选择的特征中挑出能够在测试集上获得最好效果的那个，加入到我们的特征子集中去，直到特征子集数目达到预设值或效果没有明显变化。
- 思考优缺点：
 - 优点：快， $O(nk)$ 的时间复杂度（ n 为特征全集长度， k 为子集最大长度），在大部分情况下可以取得和复杂算法相似的效果
 - 缺点：没有理论的效果保证，最差的情况表现糟糕
 - 弥补措施：多次执行取交集/ 特征分组，以组为粒度选择，增加稳定性/ 加入交叉验证

3/29/2021



特征子集评价

57

□ 如何评价特征子集？

不同的特征选择算法不仅对特征子集评价标准不同，有的还需要结合后续的学习算法模型。因此根据特征选择中子集评价标准和后续算法的结合方式主要分为过滤式(Filter)、封装式(Wrapper)和嵌入式(Embedded)三种。

1. 过滤式(Filter)评价策略方法

Filter 方法是一种计算效率较高的方法，它独立于后续的学习算法模型来分析数据集的固有的属性。通过采用一些基于信息统计的启发式准则来评价特征子集。启发式的评价函数主要分为四类：距离度量、信息度量、依赖性度量以及一致性度量。

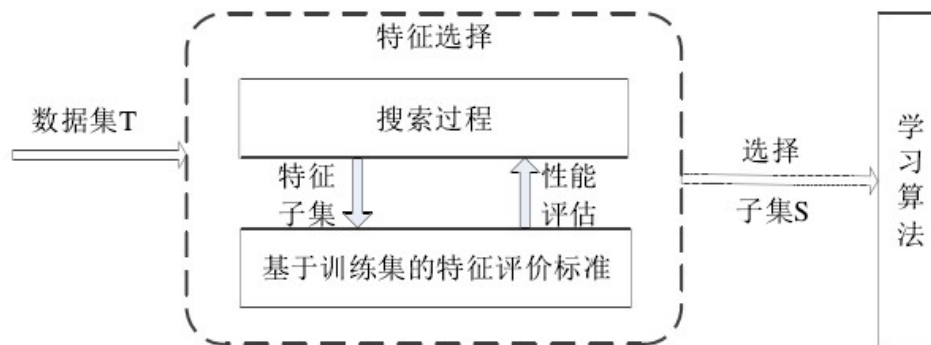


图 2-2 Filter 特征选择

从数据集本身的内在性质获得，与特定的算法无关，因此具有较好通用性，可适用大规模数据集，快速去除大量不相关特征。但倾向于选择冗余的变量，选择出的特征子集的规模也会比较大，因此常被用为特征选择的预处理方法。



特征子集评价

58

□ 如何评价特征子集？

2. 封装式(Wrapper)评价策略方法

Wrapper选择算法将特征选择作为学习算法一个组成部分，需要结合后续的学习算法，并直接将学习算法的分类性能作为特征重要性的评价标准。Wrapper 选择方法直接使用**分类器的性能作为评价的标准**，选出来的特征子集对分类一定有最好的性能。

相对于Filter 选择方法，Wrapper 方法所选择的特征子集的规模要小得多，有利于关键特征的辨识，模型的**分类性能更好**。但Wrapper 方法泛化能力较差，当改变学习算法时，需要针对该学习算法重新进行特征选择，算法的计算**复杂度高**。

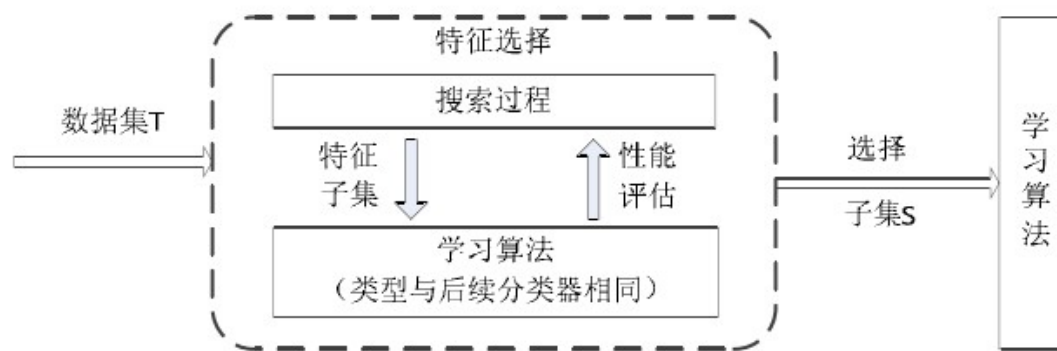


图 2-3 Wrapper 特征选择



特征子集评价

59

□ 如何评价特征子集？

3. 嵌入式(Embedded)评价策略方法

基于Embedded 嵌入式特征选择方法结合了学习算法和特征选择机制去评价学习过程中被考虑的特征。特征选择算法嵌入到学习和分类算法中，也就是**特征选择是算法模型中的一部分**，算法模型训练和特征选择**同时进行，互相结合**（即，算法具有自动进行特征选择的功能）。常见的方法有：

1). 带惩罚项的特征选择方法

其基本思想就是在模型损失函数上加上一个惩罚项，模型训练时通过惩罚项来**对特征的系数进行惩罚处理**，而在特征选择方法中常使用的是L1 正则化(regularization)项。

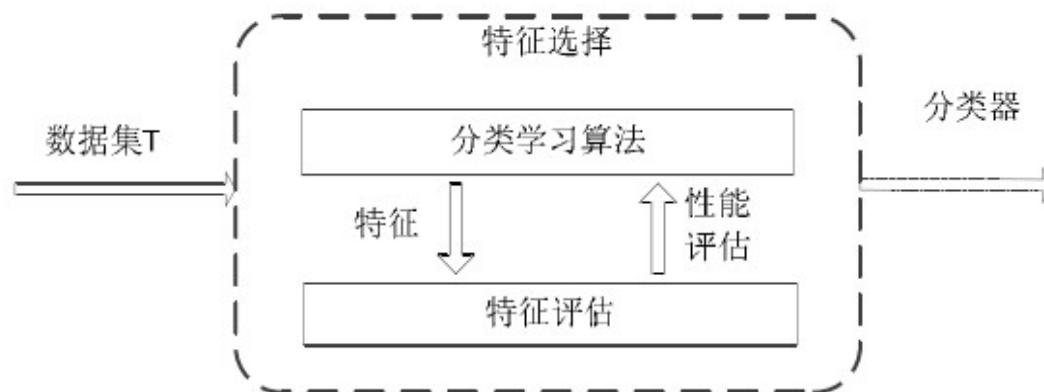


图 2-4 Embedded 特征选择



特征子集评价

60

□ 如何评价特征子集？

3. 嵌入式(Embedded)评价策略方法

基于Embedded 嵌入式特征选择方法结合了学习算法和特征选择机制去评价学习过程中被考虑的特征。特征选择算法嵌入到学习和分类算法中，也就是**特征选择是算法模型中的一部分**，算法模型训练和特征选择**同时进行，互相结合**（即，算法具有自动进行特征选择的功能）。常见的方法有：

1). 带惩罚项的特征选择方法

其基本思想就是在模型损失函数上加上一个惩罚项，模型训练时通过惩罚项来**对特征的系数进行惩罚处理**，而在特征选择方法中常使用的是L1 正则化(regularization)项。

正则化是把额外的约束或者惩罚项加到已有模型（损失函数），以防止过拟合并提高泛化能力。

损失函数由原来的 **$E(X,Y)$** 变为 **$E(X,Y)+\alpha\|w\|_1$** ,

w 是模型系数组成的向量（有些地方也叫参数parameter, coefficients）， $\|\cdot\|$ 一般是L1或者L2范数， α 是一个可调的参数，控制着正则化的强度。当用在线性模型上时，L1正则化和L2正则化也称为Lasso和Ridge。



特征子集评价

61

□ 如何评价特征子集？

2). 基于树模型的特征选择方法

这些算法在树增长过程的每一步都必须选择一个特征，将样本集划分为纯度更高的子集，而每次选择出的都是使划分效果最佳的特征，所以**决策树的生成过程就是特征选择的过程**。当决策树完全生成后，每个结点分裂所使用的特征组成的集合就是最后筛选出的特征子集。比如在比赛中经常使用的**迭代决策树(GBDT)**、**随机森林(RF)**等算法。

□ 举例：

- 前面初步筛选得到的200维特征，将其输入xgboost(一种高效的梯度提升机（GBM, Gradient boosting machine）算法)
- 训练得到特征重要性，也就是分裂树节点时起到的作用权重，自行划分阈值选取特征子集
- 为了保证不遗漏重要特征，这里不妨将树的深度设高一些

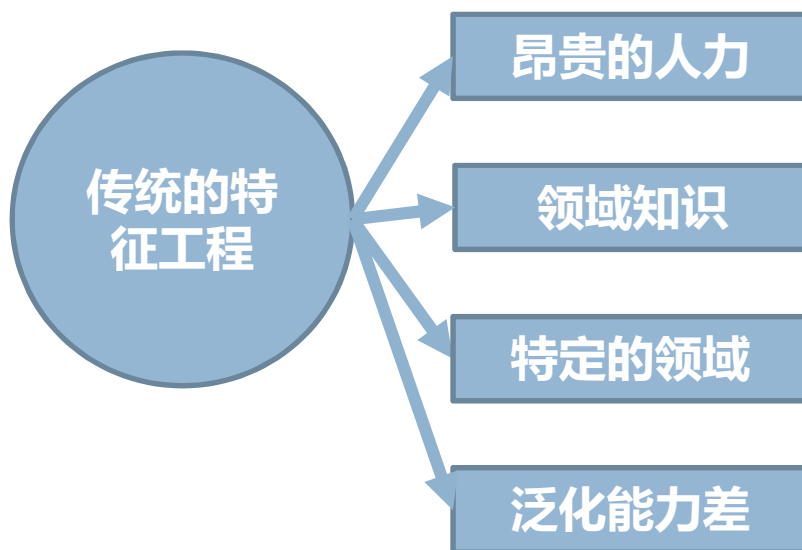
3/29/2021



传统特征工程的缺点

62

□ 传统特征工程的缺点

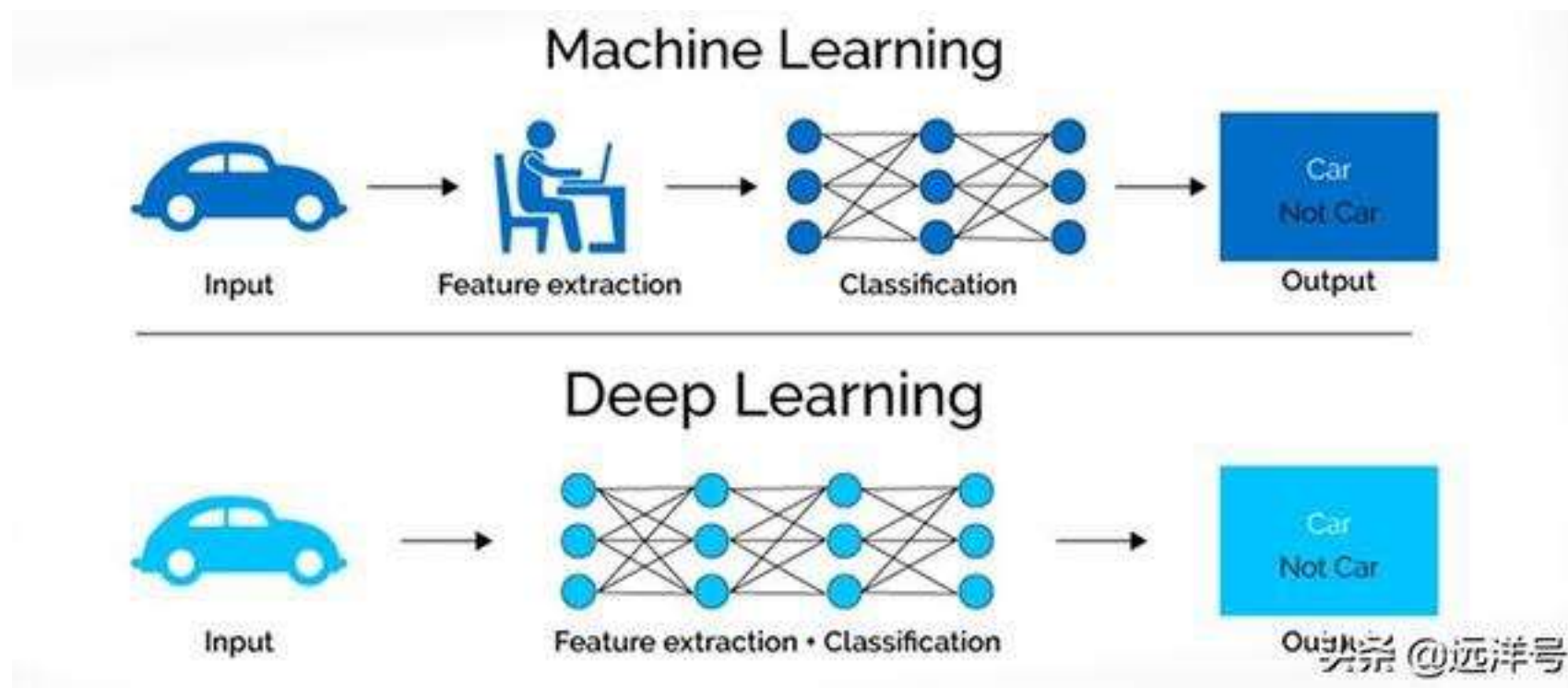




传统特征工程的缺点

63

□ 传统特征工程的缺点





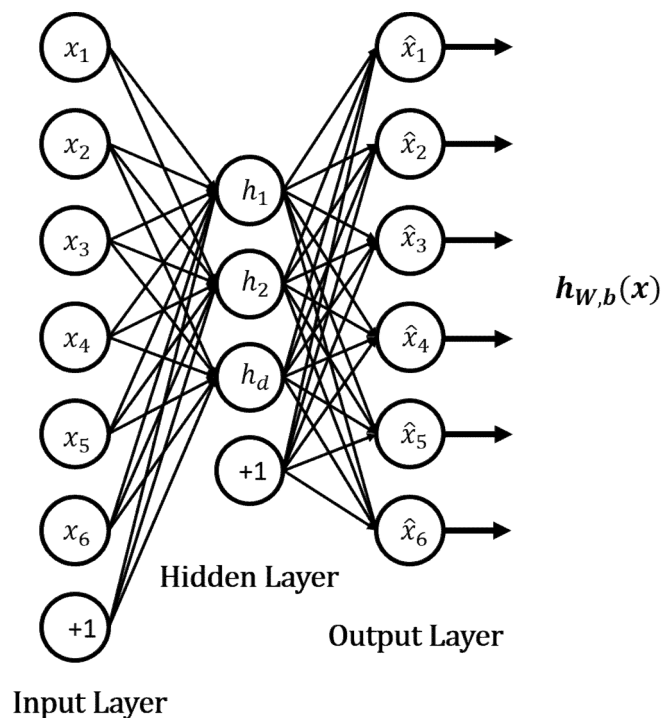
特征学习

64

□ 特征学习

如何从数据中能够自主的学习特征，在这里我们主要介绍在深度学习中常用的三种网络结构。

□ 自编码结构(Auto-Encoder)



将数据的特征 X 作为Input Layer输入
同样将原始数据特征 X' 作为Output Layer的输出来重构出原数据。

$$\text{Encoder: } H = f(A * X + b)$$

$$\text{Decoder: } X' = f(A' * H + b')$$

将中间的隐含层 H 的输出作为学习到的数据特征。

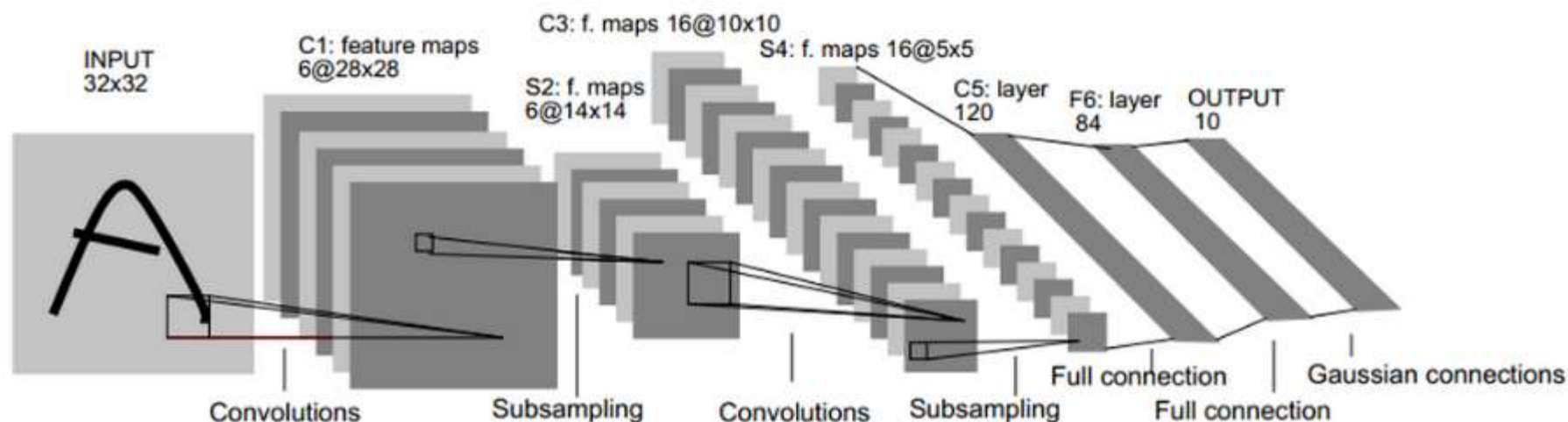
3/29/2021



特征学习

65

- 卷积神经网络(CNN): 常用于图像特征提取



卷积层：通过局部平移，利用不同的卷积核来提取图像中不同的特征

池化层：计算某个区域的特征，提高模型的泛化能力

全连接层：通过多层的神经网络，抽取更高阶的特征。

最终**全连接层的输出**即为该图像的特征向量表示。

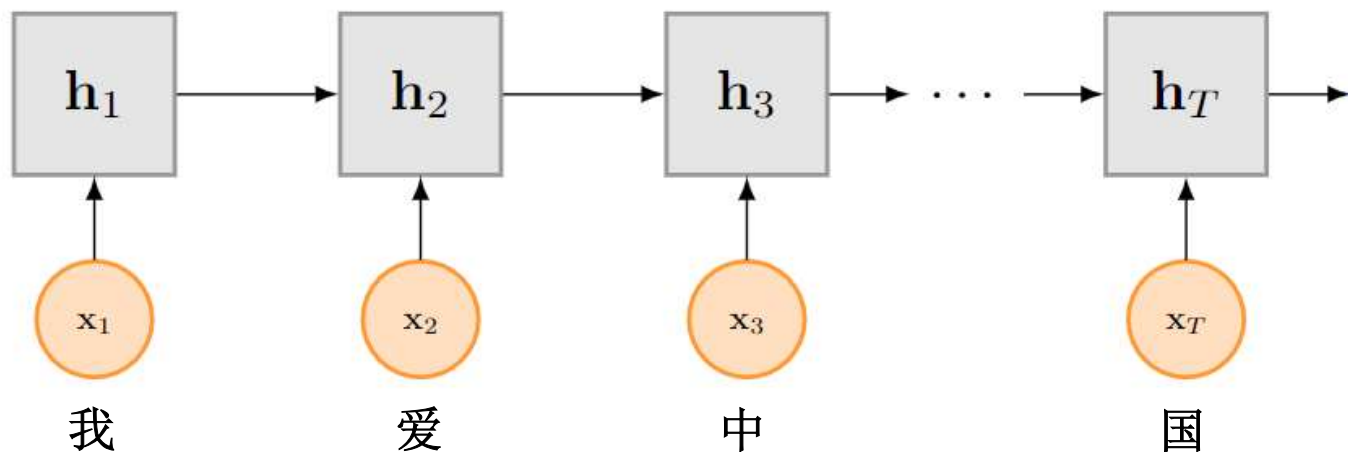
3/29/2021



特征学习

66

- 循环神经网络(RNN): 常用于序列数据的特征提取



将序列中的每个数据依次作为RNN的输入，如上图中的文本数据‘我’、‘爱’、‘中’、‘国’，并将最后一层网络的输出 h_T 作为最终序列数据的特征向量



特征学习

67

- 利用标准数据集进行特征学习（预训练）
 - 作用：模型效果验证 & 应用问题中的模型预训练
 - 图像数据预训练：ImageNet
 - <http://www.image-net.org/>
 - 1400万张图片数据，2万类别，已标注
 - 常用模型：ResNet, AlexNet, VGG等
 - 常见应用：图像分类、目标检测、目标定位
 - 文本数据预训练：Twitter, Wiki
 - <https://nlp.stanford.edu/projects/glove/>
 - 2 Billion tweets, 27 Billion 词数, 1.2M 词表
 - 常用模型：CBOW, Skip-gram, Glove等Word2vec模型
 - 常见应用：文本分类，文本推理，翻译等

训练好的特征即可直接作为其它模型的输入来使用

3/29/2021



课外实践：案例学习

68

□ IRIS(鸢尾花) + sklearn特征工程案例

□ <http://www.cnblogs.com/jasonfreak/p/5448385.html>

□ 1. 数据集的描述与导入

数据的特征:

花萼长度

花萼宽度

花瓣长度

花瓣宽度

花的类别:

山鸢尾

杂色鸢尾

维吉尼亚鸢尾



```
1 from sklearn.datasets import load_iris
2
3 #导入数据集IRIS
4 iris = load_iris()
5
6 #特征矩阵
7 iris.data
8
9 #目标向量
10 iris.target
```



课外实践：案例学习

69

□ □ 2. 数据集的预处理

a). 数据的标准化

$$x' = \frac{x - \bar{X}}{S}$$

其中 \bar{X} 为均值， S 为标准差

```
1 from sklearn.preprocessing import StandardScaler
2
3 #标准化，返回值为标准化后的数据
4 StandardScaler().fit_transform(iris.data)
```



课外实践：案例学习

70

□ □ 2. 数据集的预处理

b). 数据的归一化(规则为L2公式如下):

$$x' = \frac{x}{\sqrt{\sum_j^m x[j]^2}}$$

对特征矩阵的行处理数据，其中m为向量的维度

```
1 from sklearn.preprocessing import Normalizer
2
3 #归一化，返回值为归一化后的数据
4 Normalizer().fit_transform(iris.data)
```



课外实践：案例学习

71

□ 3. 特征的选择

a). Filter(过滤式方法)

使用**方差**做为Filter的特征评价函数，先要计算各个特征的方差，然后根据阈值选择方差大于阈值的特征。使用 sklearn 通过**方差选择法**来选择特征的代码如下：

```
1 from sklearn.feature_selection import
    VarianceThreshold
2
3 #方差选择法，返回值为特征选择后的数据
4 #参数为方差的阈值threshold
5 VarianceThreshold(threshold=3).fit_transform(iris.
    data)
```



课外实践：案例学习

72

□ 3. 特征的选择

b). Wrapper(封装式方法)

在这里我们使用递归特征消除法对一个基模型来进行多轮训练，每轮训练后，消除若干权值系数特征，再基于新的特征集进行下一轮训练，使用sklearn通过递归特征消除法来选择特征的代码如下：

```
1 from sklearn.feature_selection import RFE
2 from sklearn.linear_model import
    LogisticRegression
3
4 #递归特征消除法，返回特征选择后的数据
5 #参数为基模型estimator
6 #参数为选择的特征个数n_features_to_select
7 RFE(estimator=LogisticRegression(),
    n_features_to_select=2).fit_transform(iris.data,
    iris.target)
```




课外实践：案例学习

73

□ 3. 特征的选择

c). Embedded(嵌入式方法)

基于树模型的特征选择法有决策树、随机森林和GBDT等方法。在这里以GBDT来选择特征为例，具体sklearn的实现代码如下所示：

```
1 from sklearn.feature_selection import
   SelectFromModel
2 from sklearn.ensemble import
   GradientBoostingClassifier
3
4 #作为基模型的特征选择GBDT
5 SelectFromModel(GradientBoostingClassifier())
   .fit_transform(iris.data, iris.target)
```




参考文献

74

□ 书籍

- 数据挖掘导论
- 机器学习

□ 论文

- 《An Introduction to Variable and Feature Selection》
- 《特征选择常用算法综述》

□ 实战经验

- Pandas官方文档
- Sklearn官方文档
- Kaggle和天池比赛论坛



第二章数据分析入门小结

75

□ 数据采集

Data Acquisition

- 信息检索
- 网络爬虫
- 数据存储

□ 数据预处理

Data Preprocessing

- 数据清洗
- 数据集成
- 数据变换
- 数据规约

□ 特征工程

Feature engineering

- 特征设计
- 特征选择

3/29/2021