

一、今日内容

- a. 线性回归算法
- b. 过拟合与欠拟合
- c. 岭回归算法
- d. 特征降维
- e. Kmeans算法
- f. 购物篮分析案例

二、复习

- a. 特征预处理：处理数字类型的特征
 - i. 包含内容： 归一化 标准化
- b. 归一化：
 - i. 概念： 通过转换函数，将数字类型的特征，转换为适合机器学习的特征
 - ii. 转换公式：
 - 1. $x' = (x - \min) / (\max - \min)$
 - 2. $x'' = x' * (mx - mi) + mi$
 - iii. api: from sklearn.preprocessing import MinMaxScaler
 - 1. fit: 制定转换数据的标准
 - 2. transform: 使用制定的转换标准，对未知数据进行转换
- c. 特征抽取
 - i. 概念： 将任意类型的数据转换为适合机器学习的数字类型的特征
 - ii. 原理： 将文本中出现的特征 词的个数作为数字类的数据统计
 - iii. 中文分词工具： jieba.cut
 - 1. pip install jieba
 - 2. jieba.cut
 - iv. api: from sklearn.feature_extraction.text import CountVectorizer
 - 1. fit: 制定转换数据的标准
 - 2. transform: 对数据进行转换处理 返回值 sparse矩阵
 - 3. get_feature_names():
 - 4. toarray()
- d. 朴素贝叶斯算法
 - i. 文章分类原理： 分别计算文章属于每一个类别的概率
 - ii. 概率
 - 1. 条件概率： $P(A|B)$
 - 2. 联合概率： $P(A,B) = P(A) * P(B)$
 - iii. 贝叶斯公式

iv. 拉普拉斯平滑系数

v. api: from sklearn.naive_bayes import MultinomialNB

1. fit

2. predict

3. score

三、线性回归算法

3.1 概念

线性回归：特征值 (x) 与 目标值 (y) 之间建立的某种函数关系。

demo: $y = kx + b$: 单因子线性回归。

3.2 线性回归通式

通用公式:
$$\boxed{h(w)} = w_1x_1 + w_2x_2 + w_3x_3 \dots + b = w^T x + b$$

目标值特征值特征值特征值

3.3 回归系数和损失的计算

真实关系: 真实房子价格 = $0.02 \times \text{中心区域的距离}$ + $0.04 \times \text{城市一氧化氮浓度}$ + $(-0.12 \times \text{自住房平均房价})$ + $0.254 \times \text{城镇犯罪率}$

x1

x2

x3

x4

线性回归关系: $h(w) = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4$

question: 如何计算线性回归的回归系数。

正规方程:

概念: 通过公式一步计算出模型损失最低的回归系数。

公式:

$$w = (X^T X)^{-1} X^T y$$

缺点: 当模型过于复杂, 有可能出现无法计算出损失的情况。

error

$$\sum_{i=1}^m (h_w(x_i) - y_i)^2$$

预测值

真实值

房子价格

真实模型

制作模型

房子面积

梯度下降学习法

概念: 根据给出的学习率和学习方向, 一步一步的找到损失最低点的回归系数

公式:

$$w_1 := w_1 - \alpha \frac{\partial \text{cost}(w_0 + w_1 x_1)}{\partial w_1}$$

学习率

学习方向

target: 找到模型损失最低的回归系数。

cost

学习方向

学习率

w1'' w1' w1 w1' --> w1''

3.4 算法理解

波士顿房价预测：

真实关系：真实房子价格 = $0.02 \times \text{中心区域的距离}$ + $0.04 \times \text{城市一氧化氮浓度}$ + $(-0.12 \times \text{自住房平均房价})$ + $0.254 \times \text{城镇犯罪率}$

随机指定关系：预测房子价格 = $0.25 \times \text{中心区域的距离}$ + $0.14 \times \text{城市一氧化氮浓度}$ + $0.42 \times \text{自住房平均房价}$ + $0.34 \times \text{城镇犯罪率}$

随机值

随机值

随机值

随机值

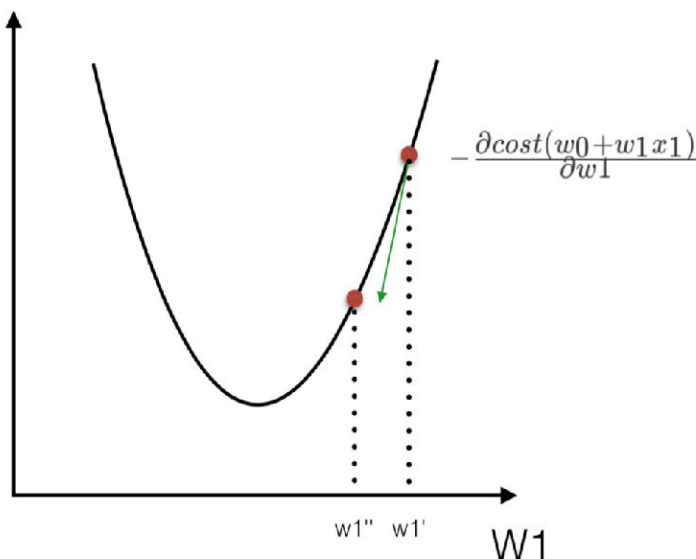
通过指定的关系制作出模型 A ,与实际情况的模型 一定有一个损失

预测的模型的损失：

$$\sum_{i=1}^m (h_w(x_i) - y_i)^2 \quad \text{cost}$$

通过调整 w 的值，使得模型的损失越来越小。

调整 w 的值，是需要 (1) 正规方程 (2) 梯度下降



3.5 梯度下降 api接口

```
1 # 导入梯度下降的api
2 from sklearn.linear_model import SGDRegressor
3 # 导入损失计算公式
4 from sklearn.metrics import mean_squared_error
5 """
6 SGDRegressor 梯度下降学习法
7
8 __init__成员函数
9
10 函数原型：
11 __init__(self, fit_intercept=True learning_rate="invscaling", eta0=0.01)
12
13 函数参数：
14     fit_intercept: bool类型,标识是否计算模型的偏置  $y = kx + b$ 
15     learning_rate: str类型: 标识学习率的计算方式。一般情况下是 constant标识固定学习率
16     eta0: 学习率
17
18 fit: 拟合训练集的数据,制作模型
19 predict: 根据数据, 进行预测结果。
20
```

```

21 mean_squared_error函数
22
23 函数原型:
24     def mean_squared_error(y_true, y_pred);
25
26 函数参数:
27     y_true: 真实值结果 y_test
28     y_Pred: 预测值的结果 测试集特征值经过模型预测出来的结果。
29
30     """

```

3.6 波士顿房价预测

实例数量:	506
属性数量:	13 数值型或类别型, 帮助预测的属性

:中位数 (第14个属性) 经常是学习目标

属性信息 (按顺序):	<ul style="list-style-type: none"> • CRIM 城镇人均犯罪率 • ZN 占地面积超过2.5万平方英尺的住宅用地比例 • INDUS 城镇非零售业务地区的比例 • CHAS 查尔斯河虚拟变量 (= 1 如果土地在河边; 否则是0) • NOX 一氧化氮浓度 (每1000万份) • RM 平均每居民房数 • AGE 在1940年之前建成的所有者占用单位的比例 • DIS 与五个波士顿就业中心的加权距离 • RAD 辐射状公路的可达性指数 • TAX 每10,000美元的全额物业税率 • PTRATIO 城镇师生比例 • B 1000(Bk - 0.63)^2 其中 Bk 是城镇的黑人比例 • LSTAT 人口中地位较低人群的百分数 • MEDV 以1000美元计算的自有住房的中位数
缺失属性值:	无
创建者:	Harrison, D. and Rubinfeld, D.L.

特征值

目标值

这是UCI ML (欧文加利福尼亚大学 机器学习库) 房价数据集的副本。 <http://archive.ics.uci.edu/ml/datasets/Housing>

该数据集是从位于卡内基梅隆大学维护的StatLib图书馆取得的。

```

1 # 导入波士顿房价的数据
2 from sklearn.datasets import load_boston
3 # 导入归一化
4 from sklearn.preprocessing import MinMaxScaler
5 # 导入线性回归
6 from sklearn.linear_model import SGDRegressor
7 # 导入分割数据集
8 from sklearn.model_selection import train_test_split
9 # 导入模型的损失计算

```

```
10 from sklearn.metrics import mean_squared_error
11
12 def main():
13     """ 预测波士顿房价 """
14
15     # 1. 获取原始数据集
16     lb = load_boston()
17
18     # 2. 确定特征值与目标值
19     x = lb.data
20     y = lb.target
21
22     # 3. 分割数据集
23     x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3)
24
25     # 4. 归一化处理数据
26     concert = MinMaxScaler()
27     # 制定转换数据的标准
28     concert.fit(x_train)
29     x_train = concert.transform(x_train)
30     x_test = concert.transform(x_test)
31
32     # 实例化 梯度下降学习法
33     estimate = SGDRegressor(fit_intercept=True, learning_rate="constant", eta0=0.01)
34
35     # 拟合数据制作模型
36     estimate.fit(x_train, y_train)
37
38     # 线性回归的模型  $y = kx + b$     k,b变量
39     coef = estimate.coef_ # coef_ 回归系数
40     intercept = estimate.intercept_ # 偏置
41
42     print("模型的回归系数: ", coef)
43     print("模型的偏置: ", intercept)
44
45     # 计算模型的损失
46     y_predict = estimate.predict(x_test)
47     error = mean_squared_error(y_test, y_predict)
48
49     print("模型的损失: ", error)
50     return 0
51
52 main()
```

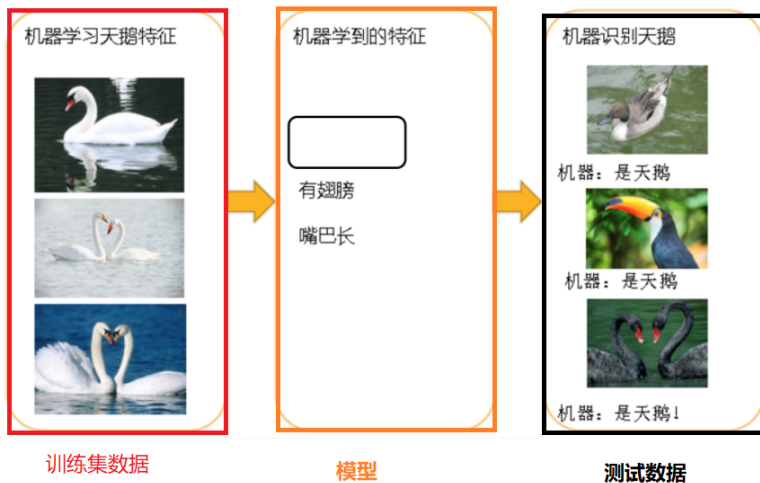
四、模型的保存与加载

```
1 # 保存模型
2 """
3 joblib.dump函数
4
5 函数原型:
6     def dump(value, filename);
7
8 函数功能:
9     将模型以文件的形式保存下来。
10
11 函数参数:
12     value: 需要保存的模型
13     filename: 模型保存的位置以及名称 模型是以 PKL结尾的
14 """
15 # joblib.dump(estimate, "./estimate.pkl")
16
17 # 导入模型
18 """
19 joblib.load函数
20
21 函数原型:
22     load(filename);
23
24 函数功能:
25     将保存的pkl文件加载到程序中
26
27 函数参数:
28     filename: 需要加载的文件名
29
30 函数返回值:
31     加载后的模型
32
33 """
34 estimate = joblib.load("./estimate.pkl")
35
```

五、过拟合与欠拟合

5.1 欠拟合

target: 制作一个识别是否是天鹅的模型!



现象:

测试集的数据进行测试的时候, 模型的准确度很低, 损失很高。

训练集的数据进行测试的时候, 模型的准确度很低, 损失很高。

这种情况就是发生了 欠拟合

原因:

a. 特征选取过少, 模型过于简单, 泛化能力太强了。

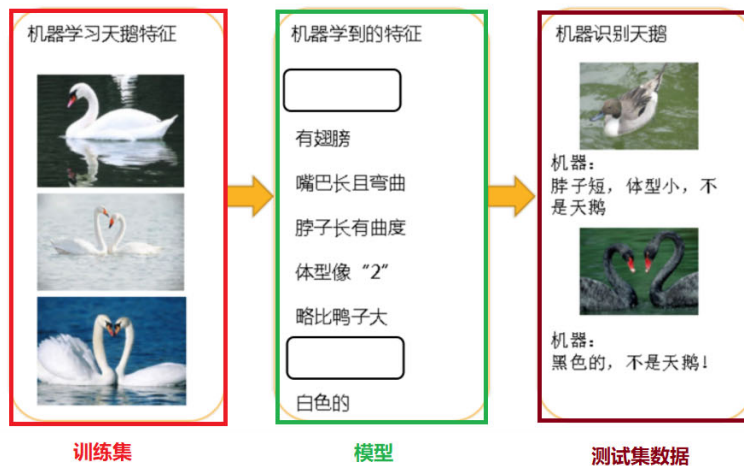
解决方式:

a. 增加数据的多样性。

b. 手动提取特征, 增加多个特征值。

5.2 过拟合

target: 制作一个识别是否是天鹅的模型



现象:

测试集数据进行测试: 模型的准确度很低, 模型的损失也很高

训练集数据进行测试: 模型的准确度很高, 模型的损失也很低。

这种情况: 发生了过拟合现象

原因:

a. 模型特征值过多, 导致模型过于复杂, 泛化能力很弱。

解决方式:

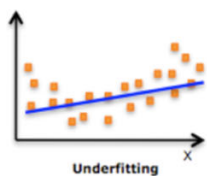
a. 手动去掉一些特征 --- 特征降维

b. 去除数据中的一些异常值。-- 噪点。

c. 正则化

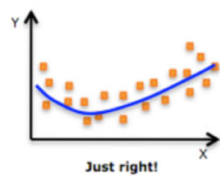
L1正则化: 将高此项系数直接置为 0

L2正则化: 将高此项系数趋近与 0



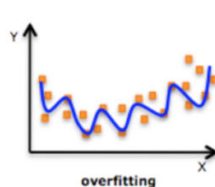
$$\theta_0 + \theta_1 x$$

欠拟合现象



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

正常拟合数据



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

过拟合现象

六、岭回归算法

6.1 概念

自带L2正则化的线性回归算法。同时可以调节正则化的力度: 正则化力度越高, 越趋近于0。

6.2 api

```
1 # 岭回归 api接口
2 from sklearn.linear_model import Ridge
```



```

3  """
4  Ridge成员函数
5
6  __init__成员函数:
7
8  函数原型:
9      def __init__(self, alpha=1.0, *, fit_intercept=True);
10
11  函数参数:
12      alpha: 正则化力度: [1.0, 10.0] alpha正则化力度越高
13      fit_intercept: 是否计算模型的偏置
14  """

```

6.3 使用岭回归制作波士顿房价模型

```

1  # 岭回归 api接口
2  from sklearn.linear_model import Ridge
3  # 导入数据
4  from sklearn.datasets import load_boston
5  # 划分数据集
6  from sklearn.model_selection import train_test_split
7  # 归一化处理数据集
8  from sklearn.preprocessing import MinMaxScaler
9
10
11
12  def main():
13
14      # 1.获取原始数据集
15      lb = load_boston()
16
17      # 2.确定特征值与目标值
18      x = lb.data
19      y = lb.target
20
21      # 3.划分数据集, 将数据集划分为训练集和测试集
22      x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3)
23
24      # 4.归一化处理数据
25      concert = MinMaxScaler()
26

```

```
27     # 制定转换数据的标准
28     concert.fit(x_train)
29
30     # 开始转换数据
31     x_train = concert.transform(x_train)
32     x_test = concert.transform(x_test)
33
34     # 5.实例化算法对象
35     estimate = Ridge(alpha=2.0)
36
37     # 6.拟合数据,制作模型
38     estimate.fit(x_train, y_train)
39
40     # 7.打印模型的回归系数和偏置
41     coef = estimate.coef_
42     intercept = estimate.intercept_
43
44     # 8.打印结果
45     print("岭回归算法的回归系数: ", coef)
46     print("岭回归算法的偏置: ", intercept)
47
48     return 0
49
50 main()
```

七、特征降维

7.1 概念

降维是指在某些限定条件下，降低随机变量(特征)个数，得到一组“不相关”主变量的过程。

降低随机变量的个数:

加载中...

不相关的主变量的过程

两个变量之间没有任何的联系

7.2 特征降维包含内容

i. 特征选择

ii. 主成分分析

7.3 特征选择

7.3.1 概念

数据中包含冗余或无关变量（或称特征、属性、指标等），旨在从原有特征中找出主要特征

加载中...

7.3.2 方差过滤法

方差：数据的混乱程度。

特征方差小：某个特征大多样本的值比较相近

特征方差大：某个特征很多样本的值都有差别

7.3.3 方差过滤法API接口

```
1 # 位置
2 from sklearn.feature_selection import VarianceThreshold
3 """
4 VarianceThreshold 特征选择类
5
6 __init__成员函数
7
8 函数功能：
9     初始化一个低方差过滤法的对象
10
11 函数原型：
12     __init__(self, threshold=0.);
13
14 函数参数：
15     threshold： 方差：过滤掉小于这个值的方差
16
```

```
17 fit: 制定转换数据的标准
18 transform:对数据进行转换处理
19
20 """
```

7.3.4 课堂练习

```
1 # 位置
2 from sklearn.feature_selection import VarianceThreshold
3 """
4 VarianceThreshold 特征选择类
5
6 __init__成员函数
7
8 函数功能:
9     初始化一个低方差过滤法的对象
10
11 函数原型:
12     __init__(self, threshold=0.);
13
14 函数参数:
15     threshold: 方差: 过滤掉小于这个值的方差
16
17 fit: 制定转换数据的标准
18 transform:对数据进行转换处理
19
20 """
21
22 # 原始数据
23 data = [[90, 2, 10, 40],
24         [60, 4, 15, 45],
25         [75, 3, 13, 46]]
26
27 # 目标去除冗余的数据
28 concert = VarianceThreshold(threshold=2)
29
30 # 制定转换数据的标准
31 concert.fit(data)
32
33 # 开始转换数据
34 result = concert.transform(data)
35
```

```
36 print("result :", result )
```

7.4 主成分分析

7.4.1 概念

高维数据转化为低维数据的过程，在此过程中可能会舍弃原有数据、创造新的变量。

高维数据转换为低维数据

.
加载中...

创造新的变量

.
加载中...

7.4.2 API接口

```
1 # API接口位置
2 from sklearn.decomposition import PCA
3 """
4 PCA 成员函数
5
6 __init__成员函数
7
8 函数原型
9     __init__(self, n_components=None);
10
11 函数参数:
12     n_components: 保留原来 x% 的信息，一般是float类型的数据。 float
13
14 fit: 制定转换数据的标准
15 transform:对数据进行转换。
16 """
```

7.4.3 购物篮分析案例

.
加载中...

i. 数据介绍

order_products__prior.csv: 订单与商品信息

i. order_id, product_id, add_to_cart_order, reordered

products.csv: 商品信息

i. product_id, product_name, aisle_id, department_id

orders.csv: 用户的订单信息

i. order_id, user_id, eval_set, order_number,

aisles.csv: 商品所属具体物品类别

i. aisle_id, aisle

i. 思路

1. 读取四张表的内容
2. 合并四张表的内容 -- 通过指定的字段
3. 将合并好的表, 纵坐标 user_id 横坐标 aisle_id 制作一张新的表

```
1 # 导入pandas模块读取四张表格
2 import pandas as pd
3
4 # 导入降维的api接口
5 from sklearn.decomposition import PCA
6
7 def main():
8
9     # 1. 读取四张表的内容 pd.read_csv()
10    aisles = pd.read_csv("./aisles.csv")
11    order_products__prior = pd.read_csv("./order_products__prior.csv")
12    orders = pd.read_csv("./orders.csv")
13    products = pd.read_csv("./products.csv")
14
15    # 2. 合并四张表里面的内容 pd.merge()
16    mt1 = pd.merge(order_products__prior, products, on=["product_id", "product_id"])
17    mt2 = pd.merge(mt1, orders, on=["order_id", "order_id"])
18
19    mt = pd.merge(mt2, aisles, on=["aisle_id", "aisle_id"])
20
21    # 3. 通过横纵坐标转换为一个新的表格 pd.crosstab()
22    data = pd.crosstab(mt["user_id"], mt["aisle"])
23
24    print("data.shape :", data.shape)
25    print(data[:5])
26
```

```
27     # 4. 实例化转换数据的对象
28     concert = PCA(n_components=0.9)
29
30     # 5.开始转换数据
31     concert.fit(data)
32     result = concert.transform(data)
33     print(result)
34     print(result.shape)
35
36
37
38
39
40     return 0
41
42
43 main()
```

八、聚类算法 kmeans

8.1 概念

kmeans算法是一个无监督学习算法(输入的数据，只有特征值，没有目标值)

.

加载中...

8.2 聚类的过程

.

加载中...

8.3 Kmeans算法api接口

```
1 # k-means算法的api接口
2 from sklearn.cluster import KMeans
3 """
4 KMeans类
5
6 __init__成员函数
7
```

```

8  函数原型:
9      __init__(self, n_clusters=8);
10
11  函数参数:
12      n_clusters: 具体聚成几个类别。
13
14  fit 成员函数
15
16  函数原型:
17      def fit(self, X);
18
19  predict成员函数
20
21  函数原型:
22      predict(self, X);
23
24
25  """

```

8.4 购物篮聚类

```

1  # 导入pandas模块读取四张表格
2  import pandas as pd
3
4  # 导入降维的api接口
5  from sklearn.decomposition import PCA
6
7  # 导入Kmeans算法
8  from sklearn.cluster import KMeans
9
10
11  def main():
12
13      # 1.读取四张表的内容  pd.read_csv()
14      aisles = pd.read_csv("./aisles.csv")
15      order_products__prior = pd.read_csv("./order_products__prior.csv")
16      orders = pd.read_csv("./orders.csv")
17      products = pd.read_csv("./products.csv")
18
19      # 2.合并四张表里面的内容  pd.merge()
20      mt1 = pd.merge(order_products__prior, products, on=["product_id", "product_id"])

```



```
21     mt2 = pd.merge(mt1, orders, on=["order_id", "order_id"])
22
23     mt = pd.merge(mt2, aisles, on=["aisle_id", "aisle_id"])
24
25     # 3. 通过横纵坐标转换为一个新的表格 pd.crosstab()
26     data = pd.crosstab(mt["user_id"], mt["aisle"])
27
28     print("data.shape :", data.shape)
29     print(data[:5])
30
31     # 4. 实例化转换数据的对象
32     concert = PCA(n_components=0.9)
33
34     # 5.开始转换数据
35     concert.fit(data)
36     result = concert.transform(data)
37
38     # 6. 实例化一个算法
39     estimate = KMeans(n_clusters=4)
40
41     estimate.fit(result[:500])
42
43     # 查看前500个数据的类别
44     temp = estimate.predict(result[:500])
45     print(temp)
46
47
48
49
50
51     return 0
52
53
54 main()
```