

机器学习

概述

什么是机器学习

机器学习是一门能够让编程计算机从数据中学习的计算机科学。

一个计算机程序在完成任务T之后，获得经验E，其表现效果为P，如果任务T的性能表现，也就是用以衡量的P，随着E增加而增加，那么这样计算机程序就被称为机器学习系统。

自我完善，自我增进，自我适应。

为什么需要机器学习

- 自动化的升级和维护
- 解决那些算法过于复杂甚至跟本就没有已知算法的问题

机器学习的问题

1. 建模问题

所谓机器学习，在形式上可这样理解：在数据对象中通过统计或推理的方法，寻找一个接受特定输入X，并给出预期输出Y的功能函数f，即 $Y=f(X)$ 。

2. 评估问题

针对已知的输入，函数给出的输出(预测值)与实际输出(目标值)之间存在一定的误差，因此需要构建一个评估体系，根据误差的大小判定函数的优劣。

3. 优化问题

学习的核心在于改善性能，通过数据对算法的反复锤炼，不断提升函数预测的准确性，直至获得能够满足实际需求的最优解，这个过程就是机器学习。

机器学习的种类

监督学习、无监督学习、半监督学习、强化学习

1. 有监督学习：用已知输出评估模型的性能。
2. 无监督学习：在没有已知输出的情况下，仅仅根据输入信息的相关性，进行类别的划分。
3. 半监督学习：先通过无监督学习划分类别，再根据人工标记通过有监督学习预测输出。
4. 强化学习：通过对不同决策结果的奖励和惩罚，使机器学习系统在经过足够长时间的训练以后，越来越倾向于给出接近期望结果的输出。

批量学习和增量学习

1. 批量学习：将学习的过程和应用的过程截然分开，用全部的训练数据训练模型，然后再在应用场景中实现预测，当预测结果不够理想时，重新回到学习过程，如此循环。
2. 增量学习：将学习的过程和应用的过程统一起来，在应用的同时以增量的方式，不断学习新的内容，边训练边预测。

基于实例的学习和基于模型的学习

1. 根据以往的经验，寻找与待预测输入最接近的样本，以其输出作为预测结果。

年龄	学历	经验	性别	月薪
25	硕士	2	女	中
20	本科	3	男	低
...
20	本科	3	男	?

2. 基于模型的学习：根据以往的经验，建立用于联系输出和输入的某种数学模型，将待预测输入代入该模型，预测其结果。

|输入|输出|

|1 | 2|

|2 | 4|

|3 | 6| $Y = 2 * X$

|... |

|9 | ?| $\rightarrow 18$

机器学习的一般过程

数据处理

1. 数据收集 （数据检索、数据挖掘、爬虫）
2. 数据清洗
3. 特征工程

机器学习

1. 选择模型 （算法）
2. 训练模型 （算法）
3. 评估模型 （工具、框架、算法知识）
4. 测试模型

业务运维

1. 应用模型
2. 维护模型

机器学习的典型应用

股价预测、推荐引擎、自然语言识别、语音识别、图像识别、人脸识别

机器学习的基本问题

1)回归(Regression)问题：根据已知的输入和输出寻找某种性能最佳的模型，将未知输出的输入代入模型，得到连续的输出。

2)分类(Classification)问题：根据已知的输入和输出寻找某种性能最佳的模型，将未知输出的输入代入模型，得到离散的输出。

3)聚类问题：根据已知输入的相似程度，将其划分为不同的群落。

4)降维问题：在性能损失尽可能小的前提下，降低数据的复杂度。

数据预处理

数据预处理的过程：输入数据 -> 模型 -> 输出数据

数据样本矩阵

年龄	学历	经验	性别	月薪
25	硕士	2	女	10000
20	本科	3	男	8000
...

一行一样本，一列一特征。

数据预处理相关库

```
1 # 解决机器学习问题的科学计算工具包
2 import sklearn.preprocessing as sp
3
4     年龄      工作经验      月薪
5 张三    23        2        10000
6 李四    21        1        9000
7 王五    40       18       10500
```

均值移除(标准化)

由于一个样本的不同特征值差异较大，不利于使用现有机器学习算法进行样本处理。均值移除可以让样本矩阵中的每一列的平均值为0，标准差为1。

如何使样本矩阵中的每一列的平均值为0呢？

```
1 例如有一列特征值表示年龄： 17, 20, 23
2 mean = (17 + 20 + 23)/3 = 20
3 a' = -3
4 b' = 0
5 c' = 3
6 完成！
```

如何使样本矩阵中的每一列的标准差为1呢？

```
1 a' = -3
2 b' = 0
3 c' = 3
4 s' = std(a', b', c')
5 [a'/s', b'/s', c'/s']
```

均值移除API:

```
1 import sklearn.preprocessing as sp
2 # scale函数用于对函数进行预处理，实现均值移除。
3 # array为原数组，返回A为均值移除后的结果。
4 A = sp.scale(array)
```

案例：

```

1 import numpy as np
2 import sklearn.preprocessing as sp
3 raw_samples = np.array([
4     [17., 100., 4000],
5     [20., 80., 5000],
6     [23., 75., 5500]])
7
8 std_samples = sp.scale(raw_samples)
9 print(std_samples)
10 print(std_samples.mean(axis=0))
11 print(std_samples.std(axis=0))

```

范围缩放

将样本矩阵中的每一列的最小值和最大值设定为相同的区间，统一各列特征值的范围。一般情况下会把特征值缩放至[0, 1]区间。

如何使一组特征值的最小值为0呢？

```

1 例如有一列特征值表示年龄： [17, 20, 23]
2 每个元素减去特征值数组所有元素的最小值即可： [0, 3, 6]

```

如何使一组特征值的最大值为1呢？

```

1 [0, 3, 6]
2 把特征值数组的每个元素除以最大值即可： [0, 1/2, 1]

```

范围缩放API:

```

1 # 创建MinMax缩放器
2 mms = sp.MinMaxScaler(feature_range=(0, 1))
3 # 调用mms对象的方法执行缩放操作，返回缩放过后的结果
4 result = mms.fit_transform(原始样本矩阵)

```

案例:

```

1 import numpy as np
2 import sklearn.preprocessing as sp
3 raw_samples = np.array([
4     [17., 100., 4000],
5     [20., 80., 5000],
6     [23., 75., 5500]])
7 print(raw_samples)
8 mms_samples = raw_samples.copy()
9 for col in mms_samples.T:
10     col_min = col.min()
11     col_max = col.max()
12     a = np.array([
13         [col_min, 1],
14         [col_max, 1]])
15     b = np.array([0, 1])
16     x = np.linalg.solve(a, b)
17     col *= x[0]
18     col += x[1]

```

```

19 print(mms_samples)
20 # 根据给定范围创建一个范围缩放器
21 mms = sp.MinMaxScaler(feature_range=(0, 1))
22 # 用范围缩放器实现特征值的范围缩放
23 mms_samples = mms.fit_transform(raw_samples)
24 print(mms_samples)

```

归一化

有些情况每个样本的每个特征值具体的值并不重要，但是每个样本特征值的占比更加重要。

	动作	爱情	科幻
老王	20	10	5
小泽	4	2	1
老祁	15	11	13

所以归一化即是用每个样本的每个特征值除以该样本各个特征值绝对值的总和。变换后的样本矩阵，每个样本的特征值绝对值之和为1。

归一化相关API：

```

1 # array 原始样本矩阵
2 # norm 范数
3 # 11 - 11范数，向量中个元素绝对值之和
4 # 12 - 12范数，向量中个元素平方之和
5 # 返回归一化预处理后的样本矩阵
6 sp.normalize(array, norm='11')

```

案例：

```

1 import numpy as np
2 import sklearn.preprocessing as sp
3 raw_samples = np.array([
4     [17., 100., 4000],
5     [20., 80., 5000],
6     [23., 75., 5500]])
7 print(raw_samples)
8 nor_samples = raw_samples.copy()
9 for row in nor_samples:
10     row /= abs(row).sum()
11 print(nor_samples)
12 print(abs(nor_samples).sum(axis=1))
13 # 归一化预处理
14 nor_samples = sp.normalize(raw_samples, norm='11')
15 print(nor_samples)
16 print(abs(nor_samples).sum(axis=1))

```

二值化

有些业务并不需要分析矩阵的详细完整数据（比如图像边缘识别只需要分析出图像边缘即可），可以根据一个事先给定的阈值，用0和1表示特征值不高于或高于阈值。二值化后的数组中每个元素非0即1，达到简化数学模型的目的。

二值化相关API:

```
1 # 给出阈值, 获取二值化器
2 bin = sp.Binarizer(threshold=阈值)
3 # 调用transform方法对原始样本矩阵进行二值化预处理操作
4 result = bin.transform(原始样本矩阵)
```

案例:

```
1 import numpy as np
2 import sklearn.preprocessing as sp
3 raw_samples = np.array([
4     [17., 100., 4000],
5     [20., 80., 5000],
6     [23., 75., 5500]])
7 print(raw_samples)
8 bin_samples = raw_samples.copy()
9 bin_samples[bin_samples <= 80] = 0
10 bin_samples[bin_samples > 80] = 1
11 print(bin_samples)
12 # 根据给定的阈值创建一个二值化器
13 bin = sp.Binarizer(threshold=80)
14 # 通过二值化器进行二值化预处理
15 bin_samples = bin.transform(raw_samples)
16 print(bin_samples)
```

独热编码(onehot)

为样本特征的每个值建立一个由一个1和若干个0组成的序列, 用该序列对所有的特征值进行编码。

```
1 两个数    三个数    四个数
2 1         3         2
3 7         5         4
4 1         8         6
5 7         3         9
6 为每一个数字进行独热编码:
7 1-10     3-100     2-1000
8 7-01     5-010     4-0100
9         8-001     6-0010
10         9-0001
11 编码完毕后得到最终经过独热编码后的样本矩阵:
12 101001000
13 010100100
14 100010010
15 011000001
16
17
18          导演      演员      上映地点    类型
19 战狼1     吴京      余男      内地        战争
20 钢铁侠1   钢哥      铁蛋      美国        科幻
21 战狼2     吴京      吴京      内地        战争
22          10       100      10         10
23          01       010      01         01
24          10       001      10         10
25
26          101001010
```

```

27         010100101
28         100011010
29
30
31         导演      演员      上映地点      类型
32 战狼1      吴京  吴京,余男      内地      战争 主旋律
33 钢铁侠1    钢哥      铁蛋      美国      科幻 特效
34 战狼2      吴京  吴京,珊珊      内地      战争 主旋律

```

独热编码相关API:

```

1  # 创建一个独热编码器
2  # sparse: 是否使用紧缩格式（稀疏矩阵）
3  # dtype: 数据类型
4  ohe = sp.OneHotEncoder(sparse=是否采用紧缩格式, dtype=数据类型)
5  # 对原始样本矩阵进行处理, 返回独热编码后的样本矩阵。
6  result = ohe.fit_transform(原始样本矩阵)

```

```

1  ohe = sp.OneHotEncoder(sparse=是否采用紧缩格式, dtype=数据类型)
2  # 对原始样本矩阵进行训练, 得到编码字典
3  encode_dict = ohe.fit(原始样本矩阵)
4  # 调用encode_dict字典的transform方法 对数据样本矩阵进行独热编码
5  result = encode_dict.transform(原始样本矩阵)

```

案例:

```

1  import numpy as np
2  import sklearn.preprocessing as sp
3  raw_samples = np.array([
4      [17., 100., 4000],
5      [20., 80., 5000],
6      [23., 75., 5500]])
7  # 创建独热编码器
8  ohe = sp.OneHotEncoder(sparse=False, dtype=int)
9  # 用独特编码器对原始样本矩阵做独热编码
10 ohe_dict = ohe.fit(raw_samples)
11 ohe_samples = ohe_dict.transform(raw_samples)
12
13 ohe_samples = ohe.fit_transform(raw_samples)
14 print(ohe_samples)

```

标签编码

根据字符串形式的特征值在特征序列中的位置, 为其指定一个数字标签, 用于提供给基于数值算法的学习模型。

标签编码相关API:

```

1  # 获取标签编码器
2  lbe = sp.LabelEncoder()
3  # 调用标签编码器的fit_transform方法训练并且为原始样本矩阵进行标签编码
4  result = lbe.fit_transform(原始样本数组)
5  # 根据标签编码的结果矩阵反查字典 得到原始数据矩阵
6  samples = lbe.inverse_transform(result)

```

案例：

```
1 import numpy as np
2 import sklearn.preprocessing as sp
3 raw_samples = np.array([
4     'audi', 'ford', 'audi', 'toyota',
5     'ford', 'bmw', 'toyota', 'ford',
6     'audi'])
7 print(raw_samples)
8 lbe = sp.LabelEncoder()
9 lbe_samples = lbe.fit_transform(raw_samples)
10 print(lbe_samples)
11 inv_samples = lbe.inverse_transform(lbe_samples)
12 print(inv_samples)
```

回归模型

线性回归

1	输入	输出
2	0.5	5.0
3	0.6	5.5
4	0.8	6.0
5	1.1	6.8
6	1.4	7.0
7	...	
8	$y = f(x)$	

预测函数： $y = w_0 + w_1 x$

x: 输入

y: 输出

w_0 和 w_1 : 模型参数

所谓模型训练，就是根据已知的 \mathbf{x} 和 \mathbf{y} ，找到最佳的模型参数 \mathbf{w}_0 和 \mathbf{w}_1 ，尽可能精确地描述出输入和输出的关系。

$$5.0 = w_0 + w_1 \times 0.5$$

$$5.5 = w_0 + w_1 \times 0.6$$

单样本误差：

根据预测函数求出输入为x时的预测值： $y' = w_0 + w_1 x$ ，单样本误差为 $\frac{1}{2}(y' - y)^2$ 。

总样本误差：

把所有单样本误差相加即是总样本误差： $\frac{1}{2}\Sigma(y' - y)^2$

损失函数：

$$\text{loss} = \frac{1}{2}\Sigma(w_0 + w_1 x - y)^2$$

所以损失函数就是总样本误差关于模型参数的函数，该函数属于三维数学模型，即需要找到一组 w_0 w_1 使得loss取极小值。

案例：画图模拟梯度下降的过程

1. 整理训练集数据，自定义梯度下降算法规则，求出 w_0 ， w_1 ，绘制回归线。


```

1 import numpy as np
2 import matplotlib.pyplot as mp
3 train_x = np.array([0.5, 0.6, 0.8, 1.1, 1.4])
4 train_y = np.array([5.0, 5.5, 6.0, 6.8, 7.0])
5 test_x = np.array([0.45, 0.55, 1.0, 1.3, 1.5])
6 test_y = np.array([4.8, 5.3, 6.4, 6.9, 7.3])
7
8 times = 1000    # 定义梯度下降次数
9 lrate = 0.01    # 记录每次梯度下降参数变化率
10 epoches = []    # 记录每次梯度下降的索引
11 w0, w1, losses = [1], [1], []
12 for i in range(1, times + 1):
13     epoches.append(i)
14     loss = (((w0[-1] + w1[-1] * train_x) - train_y) ** 2).sum() / 2
15     losses.append(loss)
16     d0 = ((w0[-1] + w1[-1] * train_x) - train_y).sum()
17     d1 = (((w0[-1] + w1[-1] * train_x) - train_y) * train_x).sum()
18     print('{:4}> w0={:.8f}, w1={:.8f}, loss={:.8f}'.format(epoches[-1],
19 w0[-1], w1[-1], losses[-1]))
20     w0.append(w0[-1] - lrate * d0)
21     w1.append(w1[-1] - lrate * d1)
22
23 pred_test_y = w0[-1] + w1[-1] * test_x
24 mp.figure('Linear Regression', facecolor='lightgray')
25 mp.title('Linear Regression', fontsize=20)
26 mp.xlabel('x', fontsize=14)
27 mp.ylabel('y', fontsize=14)
28 mp.tick_params(labelsize=10)
29 mp.grid(linestyle=':')
30 mp.scatter(train_x, train_y, marker='s', c='dodgerblue', alpha=0.5, s=80,
31 label='Training')
32 mp.scatter(test_x, test_y, marker='D', c='orangered', alpha=0.5, s=60,
33 label='Testing')
34 mp.scatter(test_x, pred_test_y, c='orangered', alpha=0.5, s=80,
35 label='Predicted')
36 mp.plot(test_x, pred_test_y, '--', c='limegreen', label='Regression',
37 linewidth=1)
38 mp.legend()
39 mp.show()

```

1. 绘制随着每次梯度下降， w_0 ， w_1 ，loss的变化曲线。

```

1 w0 = w0[:-1]
2 w1 = w1[:-1]
3
4 mp.figure('Training Progress', facecolor='lightgray')
5 mp.subplot(311)
6 mp.title('Training Progress', fontsize=20)
7 mp.ylabel('w0', fontsize=14)
8 mp.gca().xaxis.set_major_locator(mp.MultipleLocator(100))
9 mp.tick_params(labelsize=10)
10 mp.grid(linestyle=':')
11 mp.plot(epoches, w0, c='dodgerblue', label='w0')
12 mp.legend()
13 mp.subplot(312)
14 mp.ylabel('w1', fontsize=14)

```

```

15 mp.gca().xaxis.set_major_locator(mp.MultipleLocator(100))
16 mp.tick_params(labelsize=10)
17 mp.grid(linestyle=':')
18 mp.plot(epochs, w1, c='limegreen', label='w1')
19 mp.legend()
20
21 mp.subplot(313)
22 mp.xlabel('epoch', fontsize=14)
23 mp.ylabel('loss', fontsize=14)
24 mp.gca().xaxis.set_major_locator(mp.MultipleLocator(100))
25 mp.tick_params(labelsize=10)
26 mp.grid(linestyle=':')
27 mp.plot(epochs, losses, c='orangered', label='loss')
28 mp.legend()

```

1. 基于三维曲面绘制梯度下降过程中的每一个点。

```

1 import mpl_toolkits.mplot3d as axes3d
2
3 grid_w0, grid_w1 = np.meshgrid(
4     np.linspace(0, 9, 500),
5     np.linspace(0, 3.5, 500))
6
7 grid_loss = np.zeros_like(grid_w0)
8 for x, y in zip(train_x, train_y):
9     grid_loss += ((grid_w0 + x*grid_w1 - y) ** 2) / 2
10
11 mp.figure('Loss Function')
12 ax = mp.gca(projection='3d')
13 mp.title('Loss Function', fontsize=20)
14 ax.set_xlabel('w0', fontsize=14)
15 ax.set_ylabel('w1', fontsize=14)
16 ax.set_zlabel('loss', fontsize=14)
17 ax.plot_surface(grid_w0, grid_w1, grid_loss, rstride=10, cstride=10,
18                 cmap='jet')
19 ax.plot(w0, w1, losses, 'o-', c='orangered', label='BGD')
20 mp.legend()

```

1. 以等高线的方式绘制梯度下降的过程。

```

1 mp.figure('Batch Gradient Descent', facecolor='lightgray')
2 mp.title('Batch Gradient Descent', fontsize=20)
3 mp.xlabel('x', fontsize=14)
4 mp.ylabel('y', fontsize=14)
5 mp.tick_params(labelsize=10)
6 mp.grid(linestyle=':')
7 mp.contourf(grid_w0, grid_w1, grid_loss, 10, cmap='jet')
8 cntr = mp.contour(grid_w0, grid_w1, grid_loss, 10,
9                  colors='black', linewidths=0.5)
10 mp.clabel(cntr, inline_spacing=0.1, fmt='%.2f',
11          fontsize=8)
12 mp.plot(w0, w1, 'o-', c='orangered', label='BGD')
13 mp.legend()
14 mp.show()
15

```

