

人工智能基础

人工智能是研究使用计算机来模拟人的某些思维过程和智能行为（如学习、推理、思考、规划等）的学科，主要包括计算机实现智能的原理、制造类似于人脑智能的计算机，使计算机能实现更高层次的应用。

简单来说，人工智能其实就是一个机器，它可以接收输入并给出输出。

中文房间

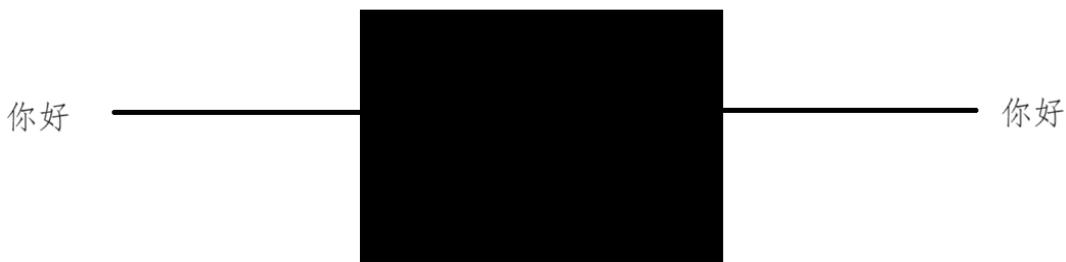
我们都知道，外国人通常认为中文是一门很难学习的语言。1980年，美国学者约翰·塞尔提出了中文房间，这有助于我们理解人工智能。

假设有一个房间，里面有一个不会说中文的外国人。房间与外界封闭，只有一个小口用于接收或发送消息。房间里提供了一个中文对照手册。另一个人在房间外向房间内发送中文，而房间里的人需要根据对照手册，向房间外的人回应。

假设它的对照手册是这样的：

输入	输出
你好	你好
你叫什么	我叫中文房间
.....

那么，如果房间外的人向房间内发送消息“你好”，房间内的人就会向房间外的人回应“你好”，如图所示：



但是，房间外的人看房间，却什么也看不到，因为房间几乎是完全封闭的。

当然，现在的人工智能通常不采用对照手册输出数据，而是通过数学方法计算得到。

训练

人工智能的训练过程就是计算方法的生成过程。训练过程中，需要给出输入和对应的期待输出，机器通过数学方法调整参数，尽可能降低实际输出与期待输出的损失，达到训练的目的。

测试

为了防止模型在训练时过于关注训练数据的特征，而无法应用在训练数据以外的数据上（通常称这种情况为过拟合），在模型训练完毕后，通常需要用不同于训练数据的数据进行测试。

paddle框架的安装

在PyCharm的下面点击Terminal，进入终端界面，依次输入以下命令：

```
1 cd .\env\Scripts\  
2 ./pip install paddlepaddle-gpu
```

普通数据处理

本节我们将以波士顿房价预测为例，展示如何使用Paddle框架进行普通数据处理。

数据集

该数据集的下载地址为archive.ics.uci.edu/ml/machine-learning-databases/housing/housing.data，打开后右击另存为，下载即可。

这个数据集共有506条数据，每个数据有13个输入和1个输出（房价），13个输入如下表所示：

字段名	类型	含义
CRIM	float	该镇的人均犯罪率
ZN	float	占地面积超过25,000平方呎的住宅用地比例
INDUS	float	非零售商业用地比例
CHAS	int	是否邻近 Charles River 1=邻近; 0=不邻近
NOX	float	一氧化氮浓度
RM	float	每栋房屋的平均客房数
AGE	float	1940年之前建成的自用单位比例
DIS	float	到波士顿5个就业中心的加权距离
RAD	int	到径向公路的可达性指数
TAX	int	全值财产税率
PTRATIO	float	学生与教师的比例
B	float	$1000(BK-0.63)^2$, 其中BK是城镇中黑人的比例
LSTAT	float	低收入人群占比

最后一个字段就是输出，也就是房价了。

接下来我们把这个数据放在项目文件夹下，先编写代码，查看以下数据情况：

```

1 import numpy as np # 用于处理数据的模块，paddle自带
2
3 dataset = np.fromfile('housing.data.txt', sep=' ') # 从文件中读取数据，sep: 分隔符
4 print(dataset)
5 print(dataset.shape) # 输出数据的形状

```

输出：

```
1 [6.320e-03 1.800e+01 2.310e+00 ... 3.969e+02 7.880e+00 1.190e+01]
2 (7084,)
```

第一行就是部分数据了，第二行指的是数据的形状，这里是一个长为7084的一维数据，也就是一条数据。但是，我们需要它是一个14列506行的二维表格，因此，我们可以使用`np.reshape`函数变换形状：

```
1 import numpy as np # 用于处理数据的模块，paddle自带
2
3 dataset = np.fromfile('housing.data.txt', sep=' ') # 从文件中读取数据，sep: 分隔符
4 dataset = dataset.reshape((506, 14))
5 print(dataset)
6 print(dataset.shape) # 输出数据的形状
```

这样它就成功变成了二维数据了，因此输出：

```
1 [[6.3200e-03 1.8000e+01 2.3100e+00 ... 3.9690e+02 4.9800e+00 2.4000e+01]
2 [2.7310e-02 0.0000e+00 7.0700e+00 ... 3.9690e+02 9.1400e+00 2.1600e+01]
3 [2.7290e-02 0.0000e+00 7.0700e+00 ... 3.9283e+02 4.0300e+00 3.4700e+01]
4 ...
5 [6.0760e-02 0.0000e+00 1.1930e+01 ... 3.9690e+02 5.6400e+00 2.3900e+01]
6 [1.0959e-01 0.0000e+00 1.1930e+01 ... 3.9345e+02 6.4800e+00 2.2000e+01]
7 [4.7410e-02 0.0000e+00 1.1930e+01 ... 3.9690e+02 7.8800e+00 1.1900e+01]]
8 (506, 14)
```

其实，`dataset`是由`numpy`模块实现的一个更快的列表，我们一般称之为数组。这样的二维数据称之为二维数组。

在`paddle`中，自定义数据集需要创建一个继承自 `paddle.io.Dataset` 类的子类，然后必须重写

`__init__`、`__getitem__` 和 `__len__` 方法：

```

1 class HouseDataset(paddle.io.Dataset):
2     def __init__(self):
3         super().__init__()
4
5     def __getitem__(self, idx):
6         pass
7
8     def __len__(self):
9         pass

```

接下来我们创建这样的子类，读取数据集：

```

1 # 文件名: dataset.py
2 import paddle
3 import numpy as np
4
5
6 class HouseDataset(paddle.io.Dataset):
7     def __init__(self, mode):
8         super().__init__()
9         dataset = np.fromfile('housing.data.txt', sep=' ', dtype=np.float32)
10        self.dataset = dataset.reshape((506, 14))
11        # 划分训练集和测试集
12        if mode == 'train':
13            self.dataset = self.dataset[: int(506*0.8)]
14        else:
15            self.dataset = self.dataset[int(506*0.8): ]
16
17    def __getitem__(self, idx):
18        return self.dataset[idx][:-1], self.dataset[idx][-1] # 返回值是元组形式，第一个
19        # 数据是特征(X)，第二个数据是预测值(Y)
20
21    def __len__(self):
22        return len(self.dataset) # 获取数据长度（数据有几条）

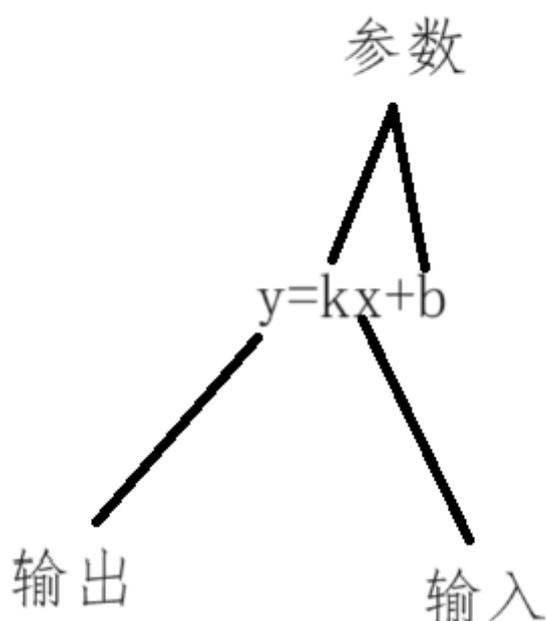
```

在数据集子类中，我们编写了划分训练集和测试集的代码，用于以后的训练和测试。另外，在读取数据时，我们设置数据类型为float32，以便paddle使用。

模型训练

模型

在paddle中，机器需要先知道计算方法，再通过训练调节参数。模型就是告诉paddle应该采用什么样的计算方法，如下图所示：



在 $y=kx+b$ 中， y 是输出， x 是输入， k 和 b 是参数，而 $y=kx+b$ 本身就是模型。

在paddle中，可以使用paddle自带的模型，也可以自定义模型。由于paddle自带模型多是关于图像处理的，因此我们采用自定义模型。

自定义模型采用paddle.nn.Sequential类，通过该类创建一个对象，这个对象就是模型：

```
1 layers = paddle.nn.Sequential(  
2     paddle.nn.Linear(13, 30), # 全连接层  
3     paddle.nn.ReLU(),  
4     paddle.nn.Linear(30, 1)  
5 )
```

这里的ReLU是激活函数，它用于防止模型线性，也就是防止模型停留在 $y=kx+b$ 上，允许模型拟合更广泛的数据。

ReLU(x)定义为：当 $x<0$ 时， $\text{ReLU}(x) = 0$ ，当 $x\geq 0$ 时， $\text{ReLU}(x) = x$ 。

训练

Paddle的训练非常简单，只需要几行代码即可实现：

```
1 import paddle
2 import dataset
3
4 ds_train = dataset.HouseDataset('train')
5 ds_test = dataset.HouseDataset('test')
6
7 # 创建模型
8 net = paddle.nn.Sequential(
9     paddle.nn.Linear(13, 30), # 全连接层
10    paddle.nn.ReLU(),
11    paddle.nn.Linear(30, 1)
12 )
13
14 # 将模型包装，以使用Paddle高级API
15 model = paddle.Model(net)
16
17 # 设置模型的优化器和损失函数
18 model.prepare(paddle.optimizer.Adam(parameters=model.parameters()), # Adam优化器
19              paddle.nn.MSELoss()) # 均方误差损失函数
20
21 # 模型训练
22 # epochs: 训练过程中完全训练多少次训练集
23 # batch_size: 训练过程中，将多少条数据同时训练
24 # verbose: 设置日志打印格式
25 model.fit(train_data=ds_train, epochs=20, batch_size=8, verbose=1)
26
27 # 模型测试
28 model.evaluate(eval_data=ds_test, verbose=1)
```

最后输出：

```
1 Epoch 1/20
2 step 51/51 [=====] - loss: 1164.6094 - 15ms/step
3 Epoch 2/20
4 step 51/51 [=====] - loss: 864.6916 - 710us/step
5 Epoch 3/20
6 step 51/51 [=====] - loss: 111.0113 - 727us/step
7 Epoch 4/20
8 step 51/51 [=====] - loss: 542.9614 - 708us/step
9 Epoch 5/20
10 step 51/51 [=====] - loss: 225.0129 - 726us/step
11 Epoch 6/20
12 step 51/51 [=====] - loss: 110.6434 - 745us/step
13 Epoch 7/20
14 step 51/51 [=====] - loss: 112.8198 - 922us/step
15 Epoch 8/20
16 step 51/51 [=====] - loss: 92.5391 - 763us/step
17 Epoch 9/20
18 step 51/51 [=====] - loss: 53.8990 - 747us/step
19 Epoch 10/20
20 step 51/51 [=====] - loss: 81.2532 - 706us/step
21 Epoch 11/20
22 step 51/51 [=====] - loss: 25.8484 - 726us/step
23 Epoch 12/20
24 step 51/51 [=====] - loss: 62.0762 - 823us/step
25 Epoch 13/20
26 step 51/51 [=====] - loss: 281.1697 - 902us/step
27 Epoch 14/20
28 step 51/51 [=====] - loss: 50.4519 - 824us/step
29 Epoch 15/20
30 step 51/51 [=====] - loss: 105.5184 - 706us/step
31 Epoch 16/20
32 step 51/51 [=====] - loss: 65.0720 - 667us/step
33 Epoch 17/20
34 step 51/51 [=====] - loss: 24.2344 - 941us/step
35 Epoch 18/20
36 step 51/51 [=====] - loss: 111.5248 - 882us/step
37 Epoch 19/20
38 step 51/51 [=====] - loss: 27.7191 - 804us/step
```

```
39 Epoch 20/20
40 step 51/51 [=====] - loss: 44.8798 - 667us/step
41 Eval begin...
42 step 102/102 [=====] - loss: 312.0678 - 382us/step
43 Eval samples: 102
```

这就是训练和测试过程的日志。

损失函数

在训练过程中，机器需要知道自己的预测值与真实值的差距，因此就有了损失函数。在回归问题中，通常使用均方误差损失函数。

回归问题就像这种问题一样，它的输出是连续的，也就是说，它可能取到24，也可能取到24加一个非常小的数，这意味着它的输出是有无限种可能的。

预测

模型的预测很简单，只需要使用 `model.predict_batch` 函数即可。需要注意的是，给这个函数传入的数据应该是一个二维数组，因为它用于预测一个批次的数据，也就是多条数据：

```
1 import numpy as np
2 pred = model.predict_batch(np.array([[0.27957, 0.00, 9.690, 0, 0.5850, 5.9260, 42.60,
3 2.3817, 6, 391.0, 19.20, 396.90, 13.59]]), dtype=np.float32))
4 print(pred[0])
```

输出：

```
1 [[28.436752]]
```

模型保存和加载

当训练好一个模型后，我们希望将其保存，需要用到 `model.save` 函数：

```
1 model.save('model')
```

加载模型时，需要用到 `model.load` 函数：

```
1 model.load('model')
```