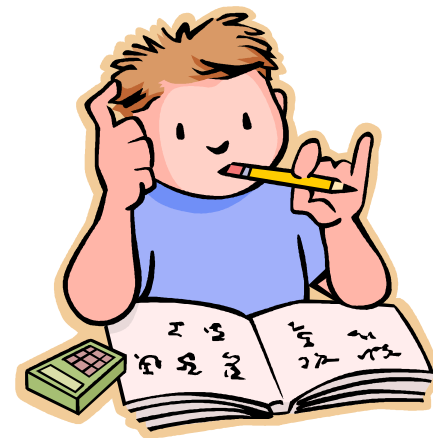




- 数据挖掘技术的产生
- 数据挖掘概念
- 数据挖掘技术的发展趋势
- 数据挖掘技术的分类问题
- 数据挖掘常用的知识表示模式与方法
- 不同数据存储形式下的数据挖掘问题
- 粗糙集方法及其在数据挖掘中的应用
- 数据挖掘的应用分析





? 大量数据、强大算力、实际需求

data



information



knowledge





- 一、C4.5: 分类决策树算法,其核心算法是ID3 算法
- 二、k-means: k-平均算法是解决聚类问题
- 三、SVM: 支持向量机
- 四、Apriori: 挖掘布尔关联规则频繁项集的算法
- 五、EM: 最大期望算法
- 六、PageRank: 网页排名、搜索引擎
- 七、AdaBoost: 自适应增强
- 八、KNN: k-近邻算法
- 九、Naive Baye: 朴素贝叶斯分类器
- 十、CART: 分类回归树



■ 网上购物时，用户推荐系统可能用到以下哪些方法？

经常一起购买的商品



- ☑ 本商品: 计算机科学丛书: 数据挖掘概念与技术(第3版) 平装 ¥61.50
- ☑ 数据挖掘导论(完整版) - 陈封能 (Pang-Ning Tan) 平装 ¥57.00
- ☑ 利用Python进行数据分析 - 麦金尼 (Wes McKinney) 平装 ¥85.50

购买此商品的顾客也同时购买

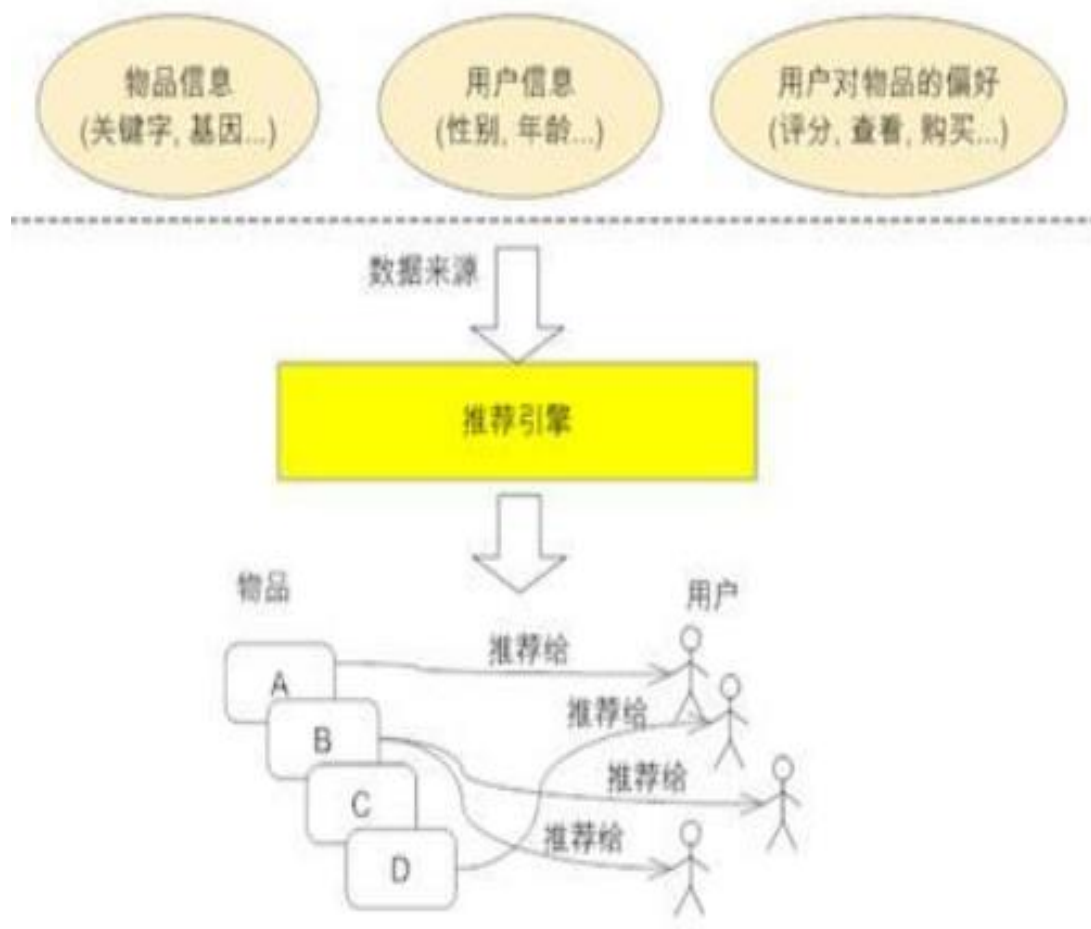


- A、决策树
- B、k-means
- C、SVM
- D、Apriori
- E、EM
- F、KNN
- G、Naive Baye
- H、其它 (填写) _____.



个性化推荐系统

是建立在海量数据挖掘基础上的一种高级商务智能平台，以帮助电子商务网站为其顾客购物提供完全个性化的决策支持和信息服务。





2、推荐系统有三个重要的模块：用户建模模块、推荐对象建模模块、推荐算法模块。

3、主要推荐方法：基于内容推荐、协同过滤推荐、基于关联规则推荐、基于效用推荐、基于知识推荐、组合推荐





发展历程：

1995年3月，卡耐基.梅隆大学提出了个性化导航系统Web Watcher；斯坦福大学推出了个性化推荐系统LIRA；

2003年，Google开创了AdWords盈利模式，通过用户搜索的关键词来提供相关的广告。AdWords的点击率很高，是Google广告收入的主要来源。

2007年，雅虎推出了SmartAds广告方案。雅虎掌握了海量的用户信息，如用户的性别、年龄、收入水平、地理位置以及生活方式等，再加上对用户搜索、浏览行为的记录，使得雅虎可以为用户呈现个性化的横幅广告。

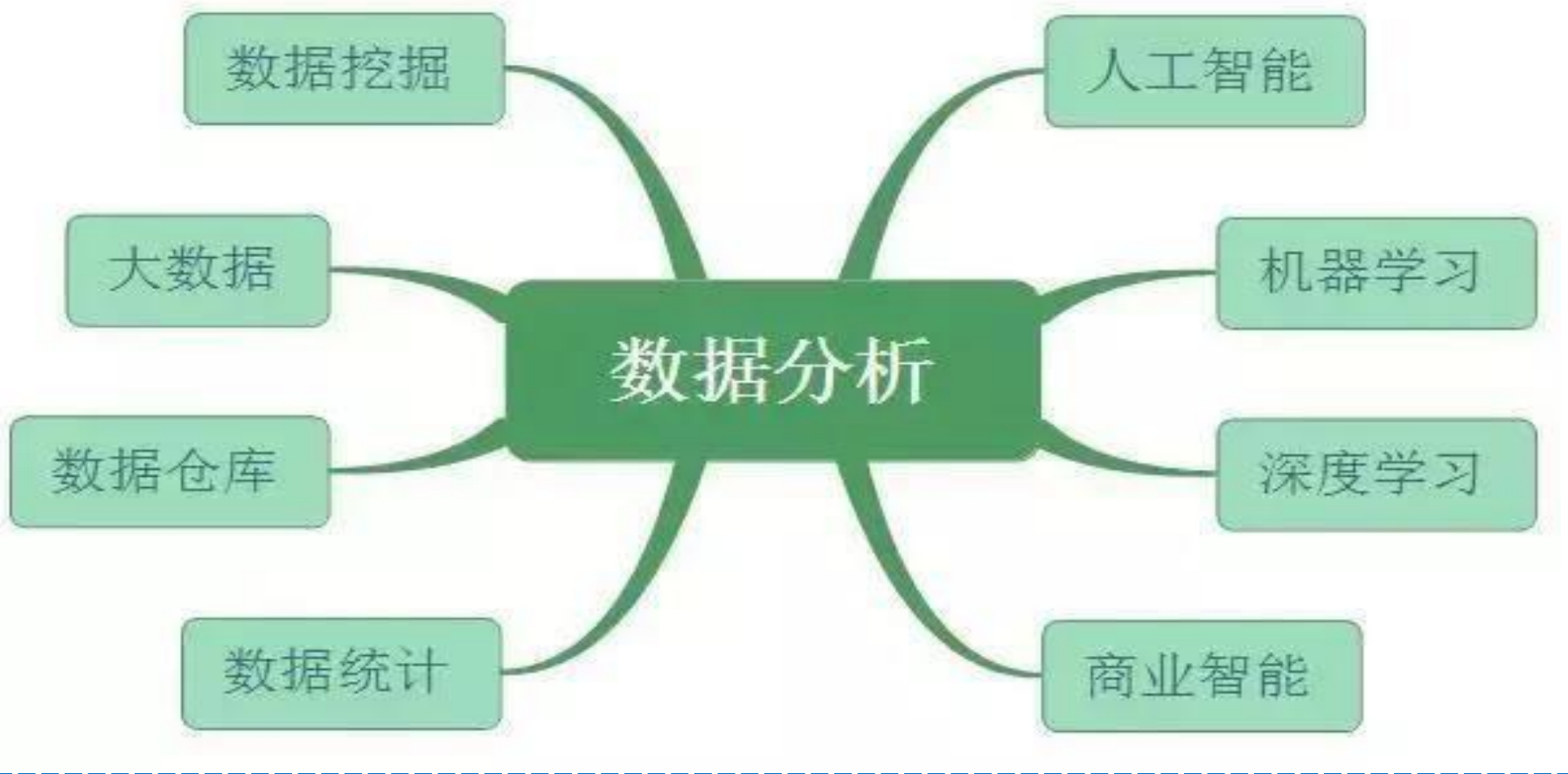


- 数据挖掘技术的产生与概念
- 数据挖掘的发展趋势
- 数据挖掘技术的分类问题
- 数据挖掘常用的知识表示模式与方法
- 不同数据存储形式下的数据挖掘问题
- 粗糙集方法及其在数据挖掘中的应用
- 数据挖掘的应用分析



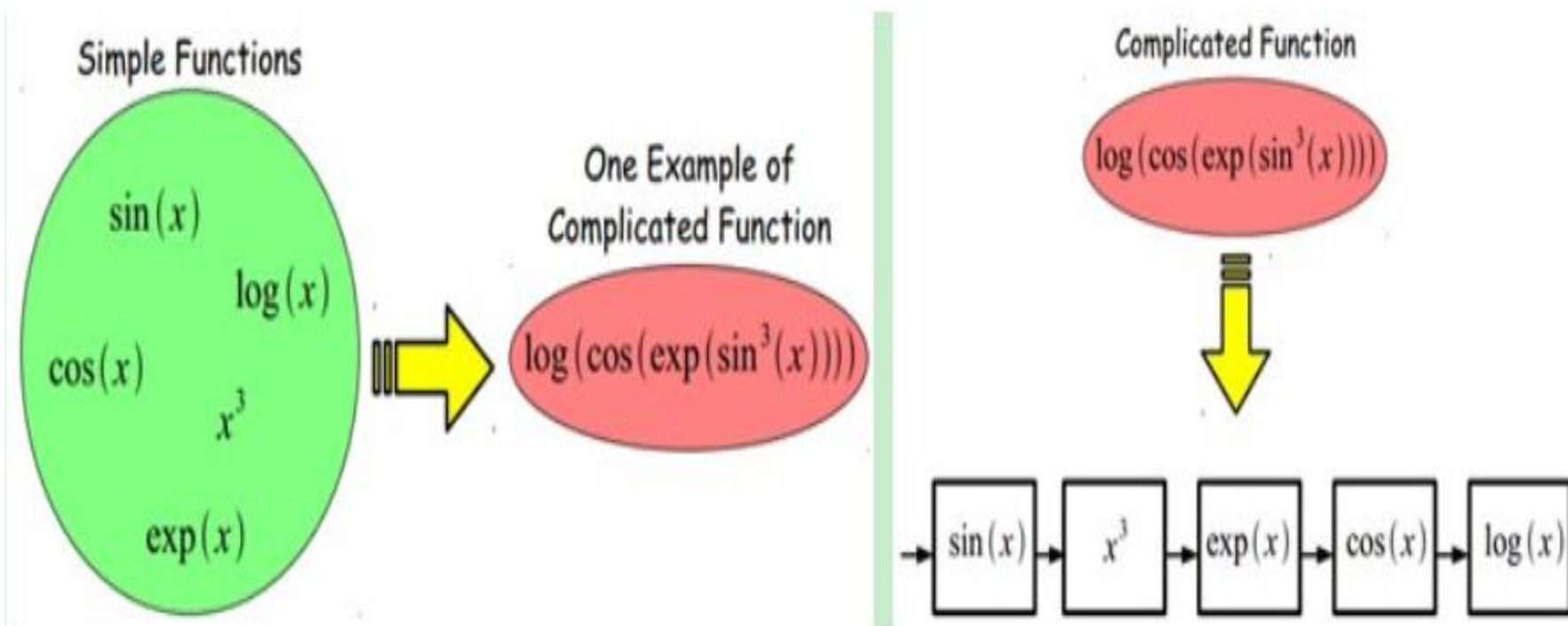


■ 数据分析是一个交叉学科





- 1、DBN: Deep Belief Network, 深度信念网络
- 2、**CNN: Convolution Neural Network**卷积神经网络
- 3、SAE: Sparse Auto Encoder, 稀疏自编码



手写数字字体识别



label = 5



label = 0



label = 4



label = 1



label = 9



label = 2



label = 1



label = 3



label = 1



label = 4



label = 3



label = 5



label = 3



label = 6



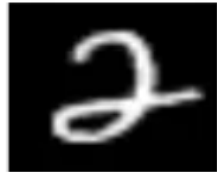
label = 1



label = 7



label = 2



label = 8



label = 6



label = 9



卷积操作



在重叠的图像和滤波器元素之间逐个进行乘法运算，按照从左向右、从上到下的顺序。

0	50	0	29
0	80	31	2
33	90	0	75
0	9	0	95

-1	0	1
-2	0	2
-1	0	1

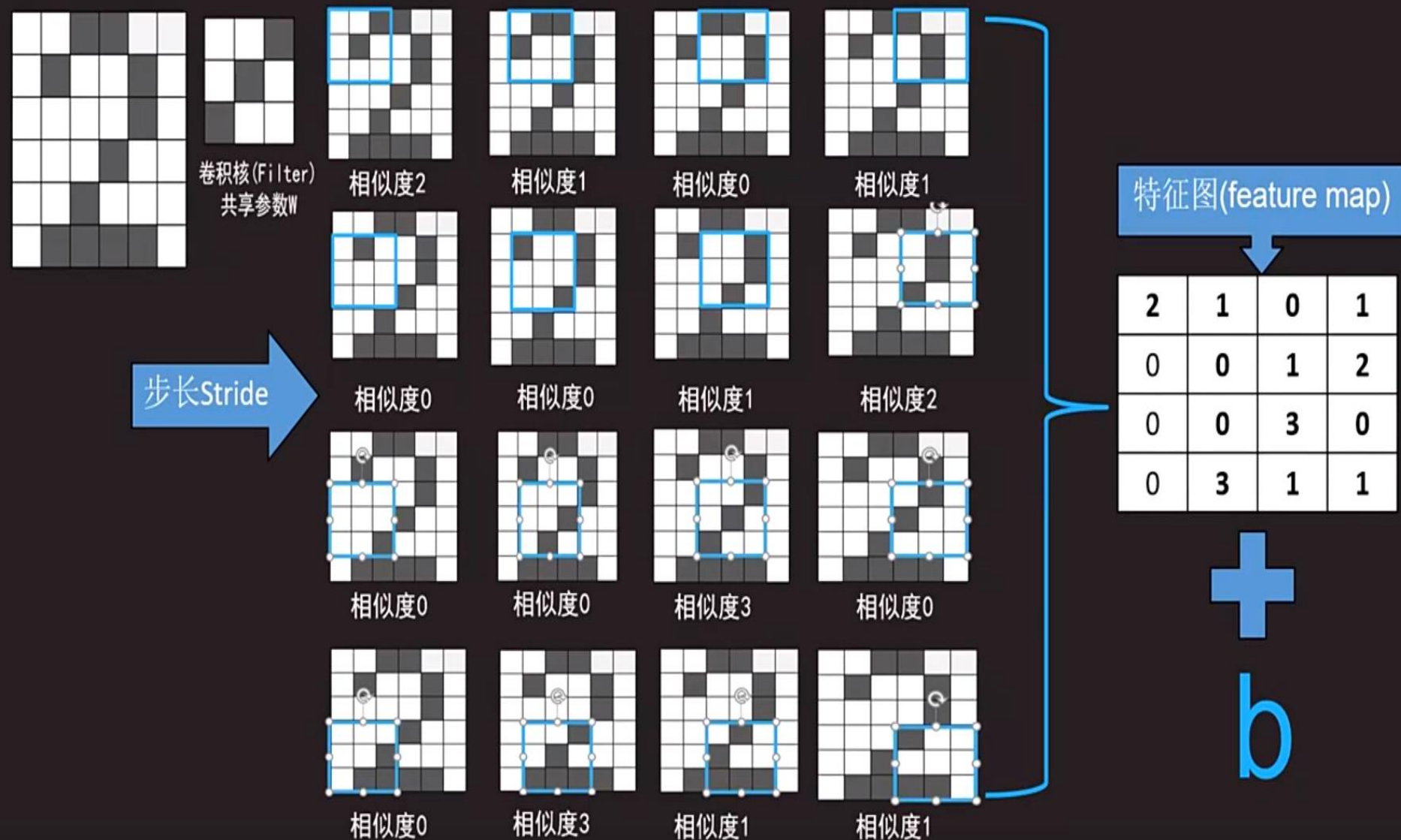
29	-192
-35	-22

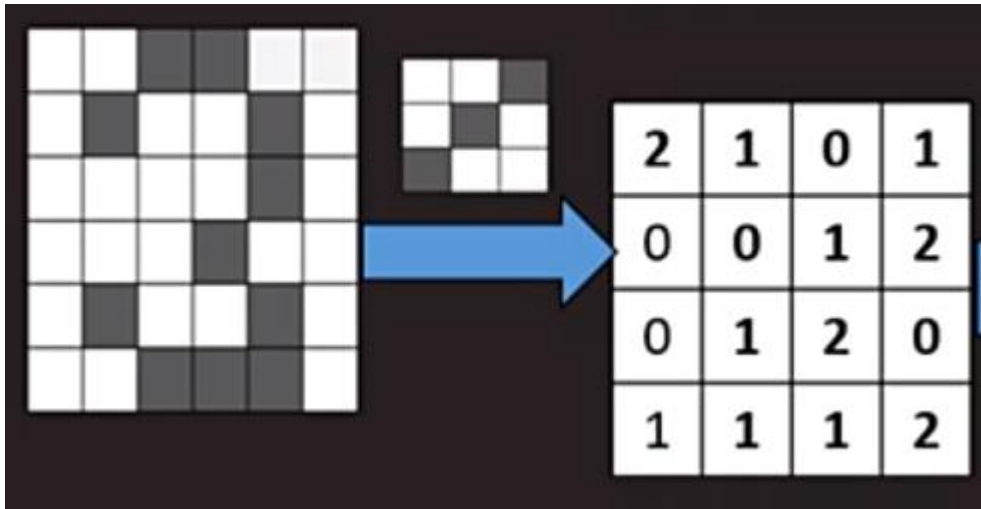
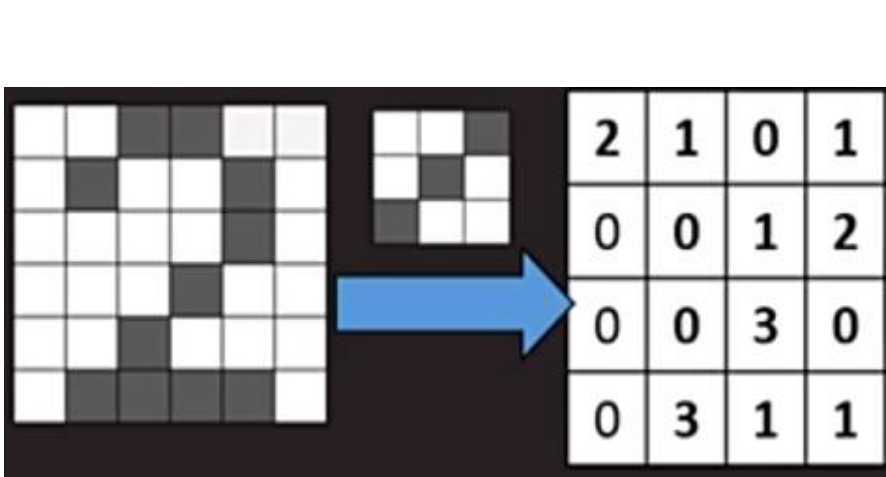
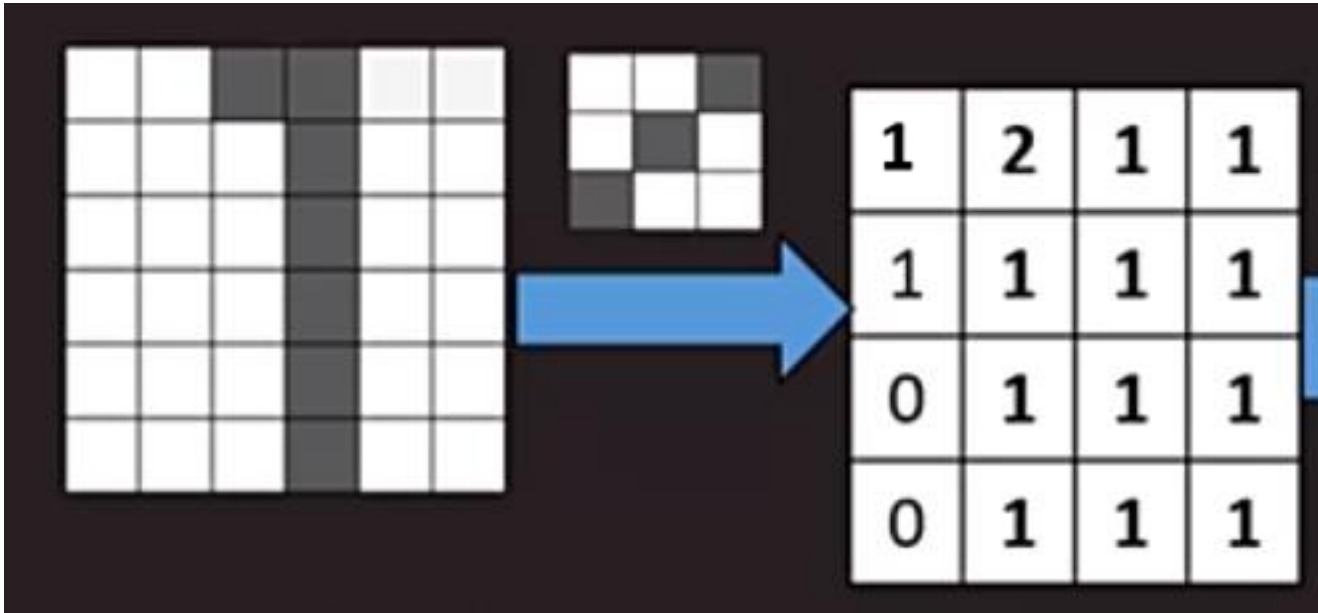


0	50	0	29
0	80	31	2
33	90	0	75
0	9	0	95

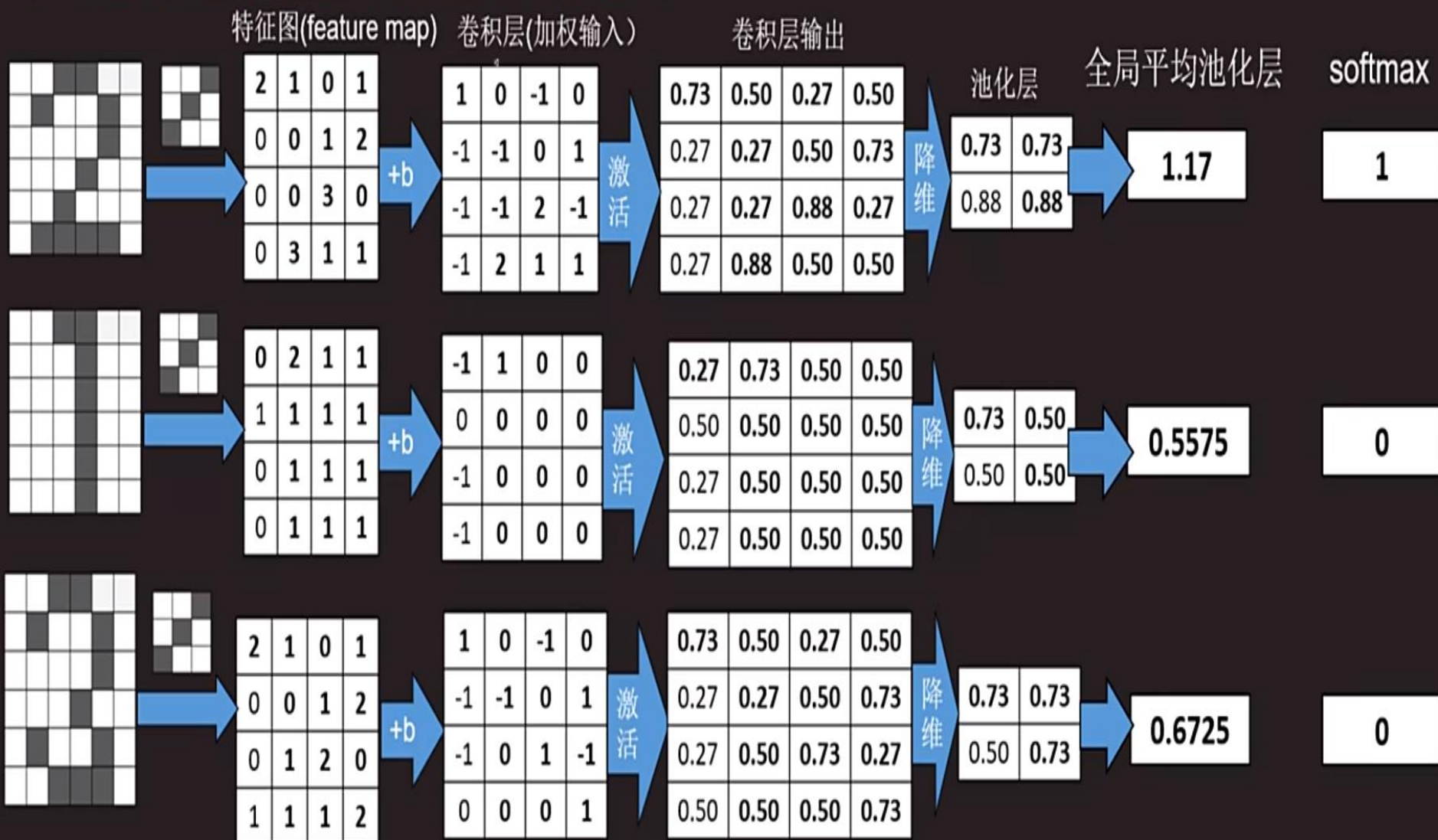
29	?
?	?

CNN——手写数字体识别





CNN——手写数字体识别





特征图(feature map)

2	1	0	1
0	0	1	2
0	0	3	0
0	3	1	1

$+b=$

卷积层(加权输入)

Z_{11}^{F1}	Z_{12}^{F1}	Z_{13}^{F1}	Z_{14}^{F1}
Z_{21}^{F1}	Z_{22}^{F1}	Z_{23}^{F1}	Z_{24}^{F1}
Z_{31}^{F1}	Z_{32}^{F1}	Z_{33}^{F1}	Z_{34}^{F1}
Z_{41}^{F1}	Z_{42}^{F1}	Z_{43}^{F1}	Z_{44}^{F1}

激活函数

α_{11}^{F1}	α_{12}^{F1}	α_{13}^{F1}	α_{14}^{F1}
α_{21}^{F1}	α_{22}^{F1}	α_{23}^{F1}	α_{24}^{F1}
α_{31}^{F1}	α_{32}^{F1}	α_{33}^{F1}	α_{34}^{F1}
α_{41}^{F1}	α_{42}^{F1}	α_{43}^{F1}	α_{44}^{F1}

池化

池化层

α_{11}^{F1}	α_{21}^{F1}
α_{32}^{F1}	α_{42}^{F1}

降维, 提取最大值

输出层识别

α_{11}^{F1}	α_{12}^{F1}	α_{13}^{F1}	α_{14}^{F1}
α_{21}^{F1}	α_{22}^{F1}	α_{23}^{F1}	α_{24}^{F1}
α_{31}^{F1}	α_{32}^{F1}	α_{33}^{F1}	α_{34}^{F1}
α_{41}^{F1}	α_{42}^{F1}	α_{43}^{F1}	α_{44}^{F1}

参数共享



特征图(feature map)

2	1	0	1
0	0	1	2
0	0	3	0
0	3	1	1

$+b=$

卷积层(加权输入)

Z_{11}^{F1}	Z_{12}^{F1}	Z_{13}^{F1}	Z_{14}^{F1}
Z_{21}^{F1}	Z_{22}^{F1}	Z_{23}^{F1}	Z_{24}^{F1}
Z_{31}^{F1}	Z_{32}^{F1}	Z_{33}^{F1}	Z_{34}^{F1}
Z_{41}^{F1}	Z_{42}^{F1}	Z_{43}^{F1}	Z_{44}^{F1}

激活函数

α_{11}^{F1}	α_{12}^{F1}	α_{13}^{F1}	α_{14}^{F1}
α_{21}^{F1}	α_{22}^{F1}	α_{23}^{F1}	α_{24}^{F1}
α_{31}^{F1}	α_{32}^{F1}	α_{33}^{F1}	α_{34}^{F1}
α_{41}^{F1}	α_{42}^{F1}	α_{43}^{F1}	α_{44}^{F1}

池化

2x2

α_{11}^{F1}	α_{12}^{F1}	α_{13}^{F1}	α_{14}^{F1}
α_{21}^{F1}	α_{22}^{F1}	α_{23}^{F1}	α_{24}^{F1}
α_{31}^{F1}	α_{32}^{F1}	α_{33}^{F1}	α_{34}^{F1}
α_{41}^{F1}	α_{42}^{F1}	α_{43}^{F1}	α_{44}^{F1}

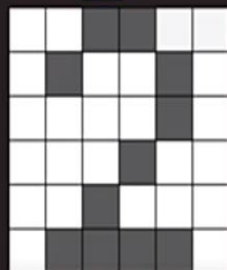
降维, 提取最大值

池化层

α_{11}^{F1}	α_{21}^{F1}
α_{32}^{F1}	α_{42}^{F1}

输出层识别

卷积核(Filter)
共享参数W



不同的卷积核，可以实现边缘检测、图像锐化、快速均值模糊、高斯模糊



Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	

求卷积有何用？



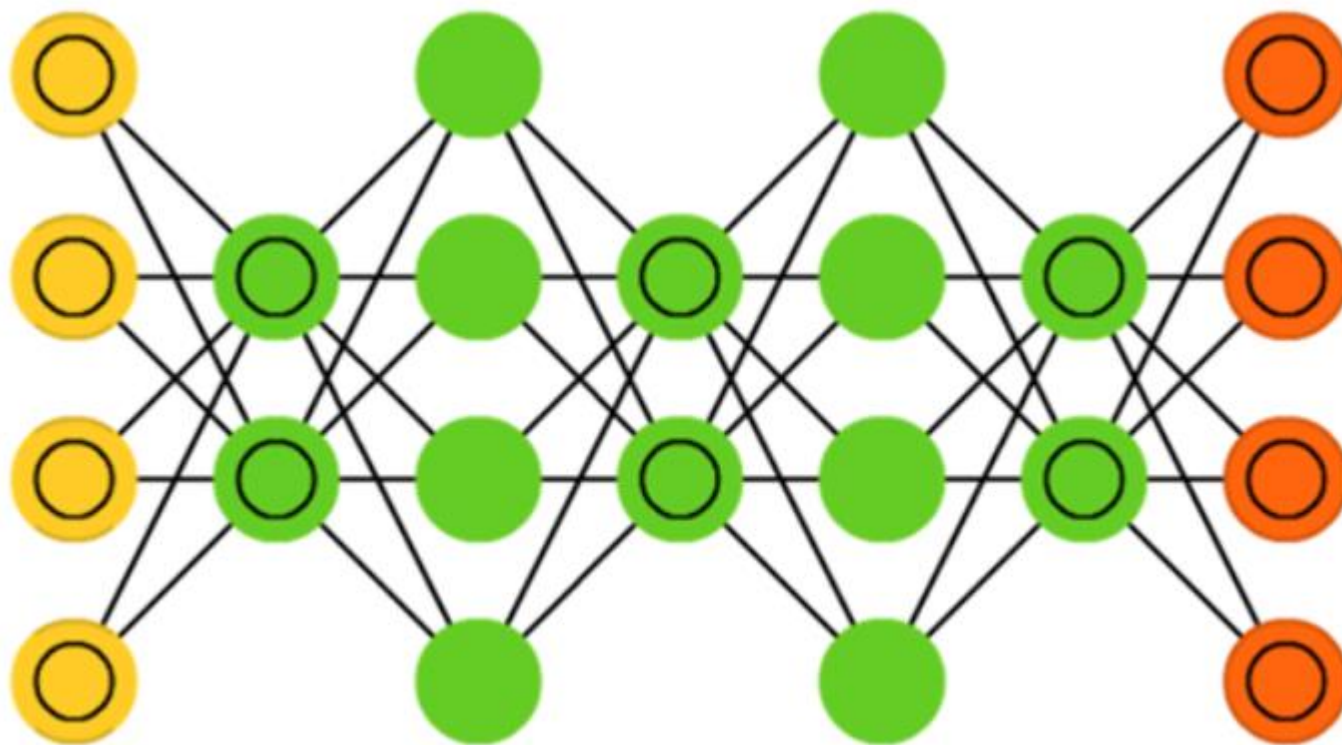
知名的Lena照片

索伯滤波器是边缘检测器。

-1	0	1
-2	0	2
-1	0	1



- 1、**DBN: Deep Belief Network**, 深度信念网络
- 2、**CNN: Convolution Neural Network**卷积神经网络
- 3、**SAE: Sparse Auto Encoder**, 稀疏自编码





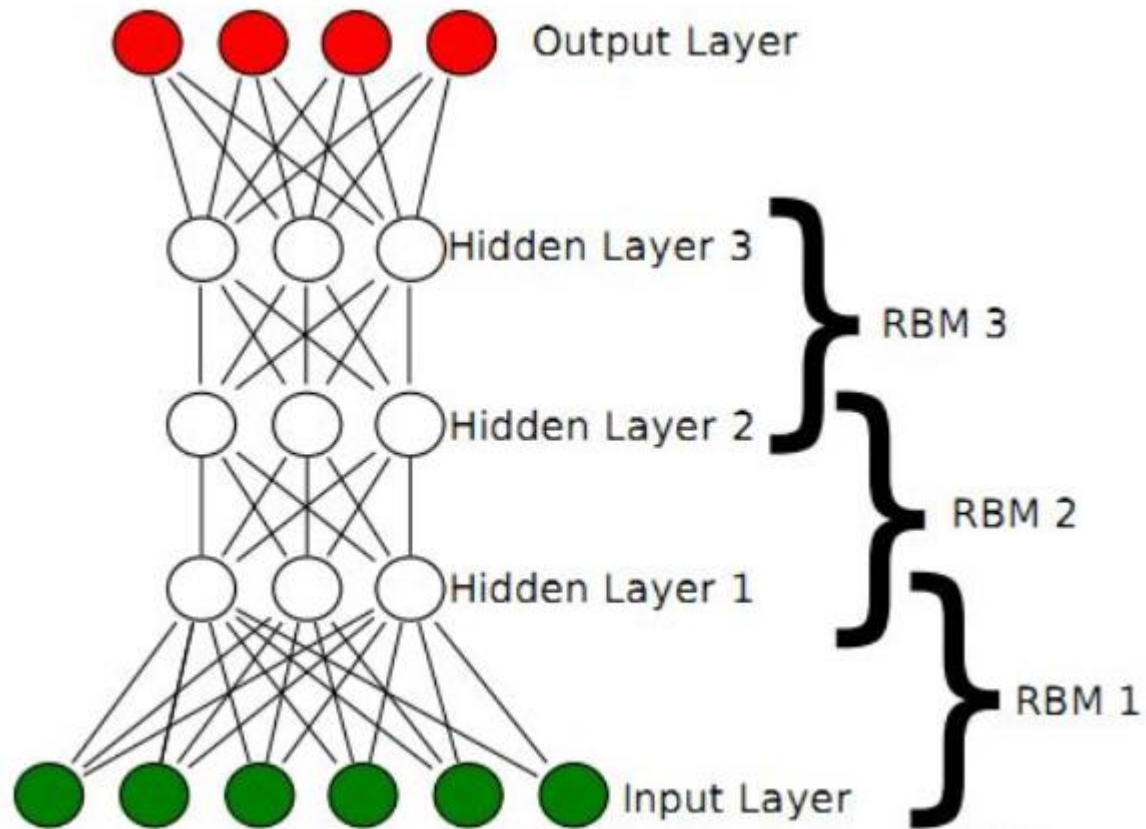
MNIST数据集中的每个图像都是28x28像素，包含一个居中的灰度数字。



MNIST训练集训练



```
2  # 28 * 28 = 784
3  input_nodes = 784
4  hidden_nodes = 200
5  output_nodes = 10
6
7  # learning rate is 0.3
8  learning_rate = 0.1
9
10 # create instance of neural network
11 n = neuralNetwork(input_nodes,hidden_nodes,output_nodes,learning_rate)
```





- 一、加载所需要的库
- 二、定义神经网络类
- 三、创建神经网络对象并用**MNIST**训练集训练
- 四、用测试集测试准确率
- 五、参数调优过程记录
- 六、测试下自己绘制的字体图片识别效果

加载所需要的库



```
1  # Code for a 3-layer neural network, and code for learning
2  # Zhouxw@ebiscn.com, 2018.8  Studying to write neural network
3  # license is GPLv2
4
5  import numpy
6  # scipy.special for the sigmoid function expit()
7  import scipy.special
8  import matplotlib.pyplot
9  # ensure the plots are inside this jupyter notebook, not
10 %matplotlib inline
11
12 # helper to load data from PNG image files
13 import imageio
14 # glob helps select multiple files using patterns
15 import glob
```

二、定义神经网络类



```
1  # neural network class definition (3 layers)
2  class neuralNetwork:
3      # initialise the neural network
4      def __init__(self,inputnodes,hiddennodes,outputnodes,learningrate):
5          # set number of nodes in each input,hidden,output layer
6          self.inodes = inputnodes
7          self.hnodes = hiddennodes
8          self.onodes = outputnodes
9          # learning rate
10         self.lr = learningrate
11
12         # link weight matrices ,wih and who
13         # weithg inside the arrays are  $w_{i,j}$ , where link is from node i to node j in the next layer
14         # w11 w21
15         # w12 w22 etc
16         self.wih = (numpy.random.normal(0.0, pow(self.hnodes,-0.5), (self.hnodes,self.inodes) ) )
17         self.who = (numpy.random.normal(0.0, pow(self.onodes,-0.5), (self.onodes,self.hnodes) ) )
18
19         # activation function is the sigmoid function
20         self.activation_function = lambda x: scipy.special.expit(x)
21
22         pass
```

```
25 def train(self,inputs_list,targets_list):
26     # convert inputs list to 2d array
27     inputs = numpy.array(inputs_list,ndmin=2).T
28     targets = numpy.array(targets_list,ndmin=2).T
29
30     # calculate signals into hidden layer
31     hidden_inputs = numpy.dot(self.wih,inputs)
32     # calculate the signals emerging from hidden layer
33     hidden_outputs = self.activation_function(hidden_inputs)
34
35     # calculate signals into final output layer
36     final_inputs = numpy.dot(self.who, hidden_outputs)
37     # calculate the signals emerging from final output layer
38     final_outputs = self.activation_function(final_inputs)
39
40     # output layer error is the (target-actual)
41     output_errors = targets - final_outputs
42     # hidden layer error is the output_errors,split by weights,recombined at hidden nodes
43     hidden_errors = numpy.dot(self.who.T, output_errors)
44
45     # update the weights for the links between the hidden and output layers
46     self.who += self.lr * numpy.dot((output_errors * final_outputs * (1.0 - final_outputs)), numpy.transpose(hidden_outputs))
47
48     # update the weights for the links between the input and hidden layers
49     self.wih += self.lr * numpy.dot((hidden_errors * hidden_outputs * (1.0 - hidden_outputs)), numpy.transpose(inputs_list))
```



```
53     # query the neural network
54     def query(self,inputs_list):
55         # convert inputs list to 2d array
56         inputs = numpy.array(inputs_list,ndmin=2).T
57
58         # calculate signals into hidden layer
59         hidden_inputs = numpy.dot(self.wih,inputs)
60         # calculate the signals emerging from hidden layer
61         hidden_outputs = self.activation_function(hidden_inputs)
62
63         # calculate signals into final output layer
64         final_inputs = numpy.dot(self.who, hidden_outputs)
65         # calculate the signals emerging from final output layer
66         final_outputs = self.activation_function(final_inputs)
67
68         return final_outputs
```


三、创建神经网络对象并用MNIST训练集训练



```
2 # 28 * 28 = 784
3 input_nodes = 784
4 hidden_nodes = 200
5 output_nodes = 10
6
7 # learning rate is 0.3
8 learning_rate = 0.1
9
10 # create instance of neural network
11 n = neuralNetwork(input_nodes,hidden_nodes,output_nodes,learning_rate)
```



```
15 # load the mnist training data csv file into a list
16 training_data_file = open("mnist_dataset/mnist_train.csv",'r')
17 training_data_list = training_data_file.readlines()
18 training_data_file.close()
19
20 # epochs is the number of times the training data set is used for training
21 epochs = 5
22 for e in range(epochs):
23     # go through all records in the training data set
24     for record in training_data_list:
25         all_values = record.split(',')
26         # scale and shift the inputs
27         inputs = (numpy.asarray(all_values[1:]) / 255.0 * 0.99) + 0.01
28         # create the target output values (all 0.01, except the desired label which is 0.99)
29         targets = numpy.zeros(output_nodes) + 0.01
```

四、用测试集测试准确率



```
4 test_data_file = open("mnist_dataset/mnist_test.csv",'r')
5 test_data_list = test_data_file.readlines()
6 test_data_file.close()
7
8 # scorecard for how well the network performs, initially empty
9 scorecard = []
10 # go through all records in the test data set
11 for record in test_data_list:
12     all_values = record.split(',')
13     # correct answer is first value
14     correct_label = int(all_values[0])
15     # scale and shift the inputs
16     inputs = (numpy.asfarray(all_values[1:]) / 255.0 * 0.99) + 0.01
17     # query the network
18     outputs = n.query(inputs)
19     # the index of the highest value corresponds to the label
20     label = numpy.argmax(outputs)
21     # print("Answer label is:", correct_label, " ; ", label, " is network's answer")
22     # append correct or incorrect to list
23     if(label == correct_label):
24         # network's answer matches correct answer, add 1 to scorecard
25         scorecard.append(1)
26     else:
27         scorecard.append(0)
28     pass
```

五、参数调优过程记录



2 有效的参数调优说明

3

4	学习率	训练轮数	隐藏层节点	结果准确率	说明
---	-----	------	-------	-------	----

5	0.3	1	100	0.9473	初始经验，效果还不错。
---	-----	---	-----	--------	-------------

6	0.6	1	100	0.9047	学习率再增加到0.6，测试准确率下降。
---	-----	---	-----	--------	---------------------

7	0.1	1	100	0.9502	降低学习率到0.1，准确率增加。
---	-----	---	-----	--------	------------------

8	0.01	1	100	0.9241	更低的学习率也不行，应该是限制了学习。
---	------	---	-----	--------	---------------------

9	0.2	1	100	0.9515	学习率调到0.2为最优
---	-----	---	-----	--------	-------------

10	0.2	5	100	0.9611	5~7轮迭代是比较好的经验值。测试准确率提高。
----	-----	---	-----	--------	-------------------------

11	0.1	5	100	0.9653	增加训练轮数，可适当降低学习率，准确率提高。
----	-----	---	-----	--------	------------------------

12	0.1	5	200	0.9723	增加隐藏层节点数，神经网络有更好的拟合能力。
----	-----	---	-----	--------	------------------------

13	0.1	5	500	0.9751	这个结果已经非常好了！
----	-----	---	-----	--------	-------------

六、测试下自己绘制的字体图片识别效果 (28*28)

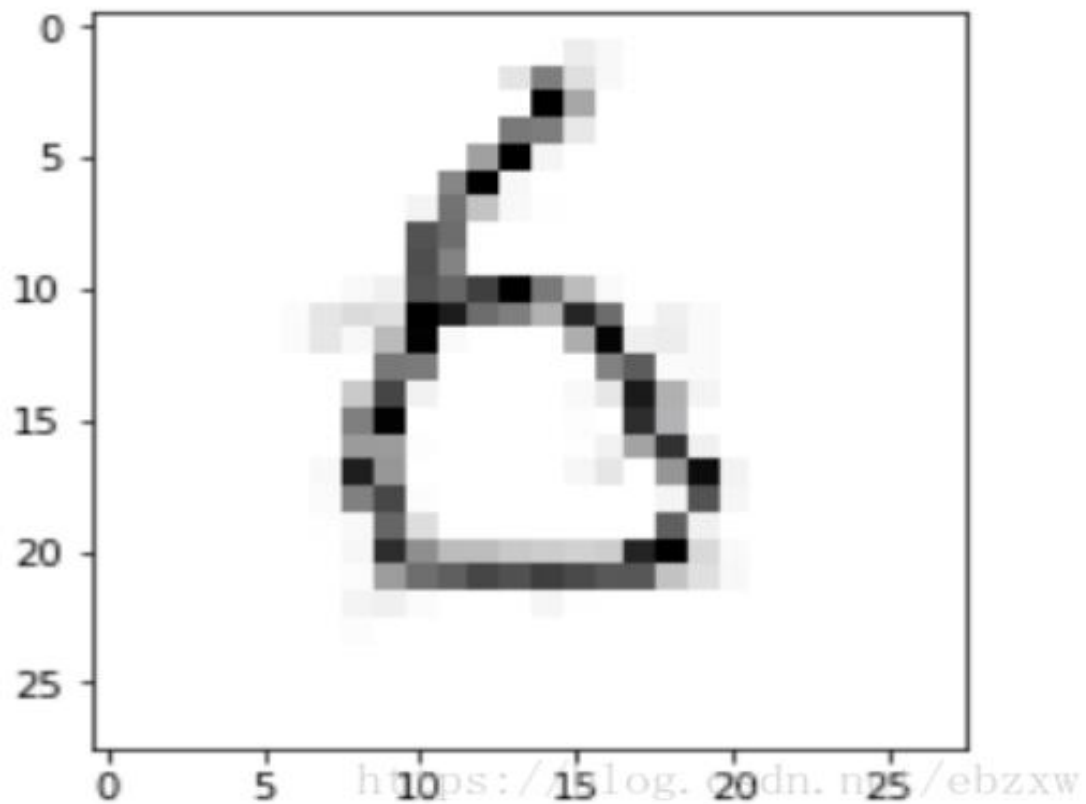


```
4  our_own_dataset = []
5
6  # Load the png image data as test data set
7  for image_file_name in glob.glob('my_own_images/2828_my_own_?.png'):
8
9      # use the filename to set the correct label
10     label = int(image_file_name[-5:-4])
11
12     # Load image data from png files into an array
13     print("loading ... ", image_file_name)
14     img_array = imageio.imread(image_file_name, as_gray=True)
15
16     # reshape from 28x28 to list of 784 values, invert values
17     img_data = 255.0 - img_array.reshape(784)
18
19     # then scale data to range from 0.01 to 1.0
20     img_data = (img_data / 255.0 * 0.99) + 0.01
21     print(numpy.min(img_data))
22     print(numpy.max(img_data))
23
24     # append label and image data to test data set
25     record = numpy.append(label, img_data)
26     our_own_dataset.append(record)
27
28     pass
```



结果样子如下:

network says 6
Good, match!





MNIST 数据集可在 <http://yann.lecun.com/exdb/mnist/> 获取, 它包含了四个部分:

- Training set images: train-images-idx3-ubyte.gz (9.9 MB, 解压后 47 MB, 包含 60,000 个样本)
- Training set labels: train-labels-idx1-ubyte.gz (29 KB, 解压后 60 KB, 包含 60,000 个标签)
- Test set images: t10k-images-idx3-ubyte.gz (1.6 MB, 解压后 7.8 MB, 包含 10,000 个样本)
- Test set labels: t10k-labels-idx1-ubyte.gz (5KB, 解压后 10 KB, 包含 10,000 个标签)

MNIST 数据集来自美国国家标准与技术研究所, National Institute of Standards and Technology (NIST). 训练集 (training set) 由来自 250 个不同人手写的数字构成, 其中 50% 是高中学生, 50% 来自人口普查局 (the Census Bureau) 的工作人员. 测试集(test set) 也是同样比例的手写数字数据.

谢谢！