

谷歌 TensorFlow 机器学习框架及应用

章敏敏 徐和平 王晓洁 周梦昀 洪淑月

(浙江师范大学 数理与信息工程学院 浙江 金华 321004)

摘要: TensorFlow 是谷歌的第二代开源的人工智能学习系统,是用来实现神经网络的内置框架学习软件库。目前,TensorFlow 机器学习已经成为了一个研究热点。由基本的机器学习算法入手,简析机器学习算法与 TensorFlow 框架,并通过在 Linux 系统下搭建环境,仿真手写字符识别的 TensorFlow 模型,实现手写字符的识别,从而实现 TensorFlow 机器学习框架的学习与应用。

关键词: TensorFlow; 机器学习; 应用

中图分类号: TP181

文献标识码: A

DOI: 10.19358/j.issn.1674-7720.2017.10.017

引用格式: 章敏敏,徐和平,王晓洁,等.谷歌 TensorFlow 机器学习框架及应用[J].微型机与应用,2017,36(10):58-60.

Application of Google TensorFlow machine learning framework

Zhang Minmin, Xu Heping, Wang Xiaojie, Zhou Mengyun, Hong Shuyue

(College of Mathematics, Physics and Information Engineering, Zhejiang Normal University, Jinhua 321004, China)

Abstract: TensorFlow is the second generation of Google's open source artificial intelligence learning system, which is used to achieve the neural network built-in framework for learning software library. At present, TensorFlow machine learning has become a hot research topic. Starting from the basic algorithm of machine learning, this paper analyzes machine learning algorithm and TensorFlow framework, and through building environment in the Linux system, simulates TensorFlow model of handwritten character recognition, to achieve handwritten character recognition, so as to realize the TensorFlow machine learning and application framework.

Key words: TensorFlow; machine learning; application

0 引言

机器学习是一门多领域交叉的学科,能够实现计算机模拟或者实现人类的学习行为,重构自己的知识结构从而改善自身的性能。2016年初,AlphaGo 以大比分战胜李世石,AI 的概念从此进入人们的视野,而机器学习就是 AI 的核心,是使计算机具有智能的根本途径。TensorFlow 是谷歌的第二代人工智能学习系统,是用来制作 AlphaGo 的一个开源的深度学习系统。

1 机器学习

可以举一个简单的例子来说明机器学习的概念,使用 k-近邻算法改进交友网站的配对效果^[1]。比如说你现在想要在交友网站上认识一个朋友,而交友网站上拥有每个注册用户的两个信息(玩视频游戏所耗时间的百分比和每年获取的飞行常客里程数),你想知道你会对哪些人比较感兴趣,这时候就可以使用机器学习算法建立一个简单的模型。可以将一些自己认为有魅力的人、魅力一般的人、不喜欢的人的这两个信息(玩视频游戏所耗时间的百分比和每年获取的飞行常客里程数)输入机器学习算法建立一个模型,如图 1 所示。当你想知道一个用户是不是你感兴趣交友的人时,输入信息,计算机通过这个模型进行计算,可以给你一个预测答案,这就是一种经典的监督

学习算法。

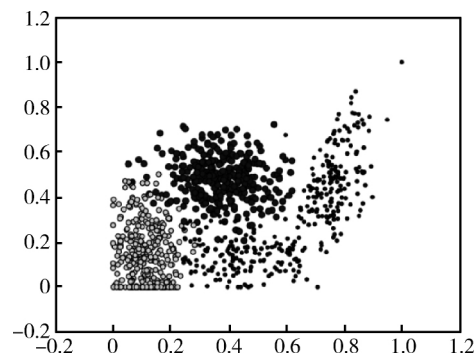


图 1 交友网站信息数据

机器学习算法有很多种类,上述例子说明的监督学习算法只是其中的一类。如果换种方式去实现这个结果,你有一堆如上的数据,但是并不对这些数据进行分类,让算法按照数据的分散方式来观察这些数据,发现数据形成了一些聚类,如图 2 所示,而通过这种方法,能够把这些数据自动地分类,这就是一种无监督学习算法。

机器学习的算法有很多,再比如用学习型算法来判断你需要多少训练信息,用什么样的更好的近似函数能够反映数据之间的关系,使得用最少的训练信息获得更准确的判断。

机器学习就是当机器想要完成一个任务,通过它不断

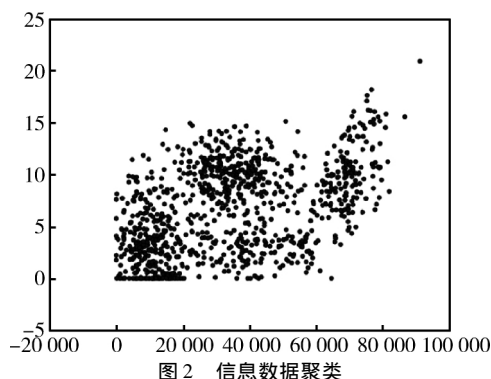


图2 信息数据聚类

地积累经验,来逐渐更好、差错减少地完成一个任务。

2 TensorFlow 的框架

2.1 TensorFlow 输入张量

TensorFlow 的命名来源于本身的运行原理。Tensor (张量) 意味着 N 维数组, Flow (流) 意味着基于数据流图的计算。用 MNIST 机器学习^[2-3] 这个例子来解释一个用于预测图片里面的数字的模型。

首先要先获得一个 MNIST 数据集,如图 3 所示,这个数据集能够在 TensorFlow 官网上进行下载。每一个 MNIST 数据单元由一张包含手写数字的图片和一个对应的标签两部分组成。把这些图片设为“xs”,把这些标签设为“ys”。MNIST 数据集拥有 60 000 行的训练数据集 (mnist.train) 和 10 000 行的测试数据集 (mnist.test)。

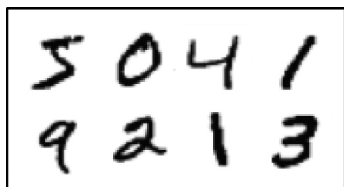


图3 数据集

每一张图片包含 28×28 个像素点。可以用一个数字数组来表示这张图片:把这个数组展开成一个向量,长度是 784。在 MNIST 训练数据集中, mnist.train.images (训练数据集中的图片) 是一个 $[60\,000, 784]$ 的张量,如图 4 所示,第一个维度数字用来对应每张图片,第二个维度数字用来索引每张图片中的像素点。在此张量里的每一个元素,都表示为某张图片里的某个像素的介于 0 和 1 之间的强度值。

相对应的标签是从 0 到 9 的数字,用来描述给定图片里表示的数字。每个数字对应着相应位置 1,如标签 0 表示为 $[1\,0\,0\,0\,0\,0\,0\,0\,0\,0]$,因此 mnist.train.labels 是

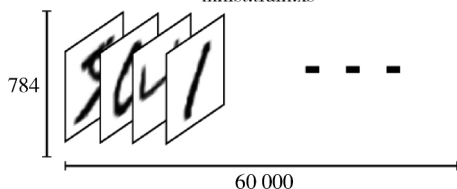


图4 image 数组

一个 $[60\,000, 10]$ 的数字矩阵,如图 5 所示。

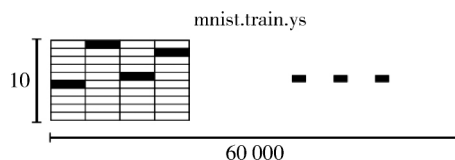


图5 标签数组

如上述的这两个数组都是二维数组,都是 TensorFlow 中的张量数据^[4],而这些数据就以流的形式进入数据运算的各个节点。而以机器算法为核心所构造的模型就是数据流动的场。TensorFlow 就是一个文件库,研究人员和计算机科学家能够借助这个文件库打造分析图像和语音等数据的系统,计算机在此类系统的帮助下,将能够自行作出决定,从而变得更加智能。

2.2 TensorFlow 代码框架

TensorFlow 是一个非常灵活的框架,它能够运行在个人计算机或者服务器的单个或多个 CPU 和 GPU 上,甚至是移动设备上。可以从上面举例的 MNIST 机器学习来分析 TensorFlow 的框架。首先,要构建一个计算的过程。MNIST 所用到的算法核心就是 softmax 回归算法,这个算法就是通过对已知训练数据同个标签的像素加权平均,来构建出每个标签在不同像素点上的权值,若是这个像素点具有有利的证据说明这张图片不属于这类,那么相应的权值为负数,相反若是这个像素拥有有利的证据支持这张图片属于这个类,那么权值是正数。

因为输入往往会带有一些无关的干扰量,于是加入一个额外的偏置量 (bias)。因此对于给定的输入图片 x 它代表的是数字 i 的证据,可以表示为:

$$\text{evidence}_i = \sum_j W_{ij} x_j + b_i \quad (1)$$

其中 W_{ij} 表示权值的矩阵, x_j 为给定图片的像素点, b_i 代表数字 i 类的偏置量。

在这里不给出详细的推导过程,但是可以得到一个计算出一个图片对应每个标签的概率大小的计算方式,可以通过如下的代码来得到一个概率分布:

$$y = \text{softmax}(Wx + b) \quad (2)$$

建立好一个算法模型之后,算法内输入的所有可操作的交互单元就像式 (2) 中的图片输入 x ,为了适应所有的图片输入,将其设置为变量占位符 placeholder。而像权重 W 和偏置值 b 这两个通过学习不断修改值的单元设置为变量 Variable。

```
train_step = tf.train.GradientDescentOptimizer(0.01).
minimize(cross_entropy)
```

TensorFlow 在这一步就是在后台给描述计算的那张图里面增添一系列新的计算操作单元来实现反向传播算法和梯度下降算法。它返回一个单一的操作,当运行这个操作时,可以用梯度下降算法来训练模型,微调变量,不

断减少成本,从而建立好一个基本模型。

建立好模型之后,创建一个会话(Session),循环1 000次,每次批处理100个数据,开始数据训练,代码如下:

```
sess = tf.InteractiveSession()
for i in range(1000):
    batch_xs, batch_ys = mnist.train.next_batch(100)
    sess.run(train_step, feed_dict={x: batch_xs, y_: batch_ys})
```

TensorFlow通过数据输入(Feeds)将张量数据输入至模型中,而张量Tensor就像数据流一样流过每个计算节点,微调变量,使得模型更加准确。

通过这个例子,可以管中窥豹了解TensorFlow的框架结构,TensorFlow对于输入的计算过程在后台描述成计算图,计算图建立好之后,创建会话Session来提交计算图,用Feed输入训练的张量数据,TensorFlow通过在后台增加计算操作单元用于训练模型,微调数据,从而完成一个机器的学习任务^[5]。

3 TensorFlow的应用

TensorFlow的支持列表里没有Windows,而人们使用的计算机大都是安装的Windows系统,虽然可以用Docker来实现在Windows上运行,但小问题很多,它支持得最好的还是基于UNIX内核的系统^[6],例如Linux,因此选择Ubuntu 15.10。

安装成功之后,可以测试一下上述MNIST_softmax的模型。在程序中加入可以判断其预测概率的代码:

```
correct_prediction = tf.equal(tf.argmax(y,1), tf.argmax(y_,1))
```

当tf.argmax(y,1)预测值与tf.argmax(y_,1)正确值相等的时候判断其为正确的预测:

```
accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
```

accuracy用来计算预测与完全错误判断之间的距离,也就是正确率,最后将它打印在显示屏上。

在导入代码之前,要先给予终端最高权限,不然在导入代码的时候会显示权限限制。成功导入代码后,命令行打印出测试结果的正确率,如图6所示为0.9191。当然

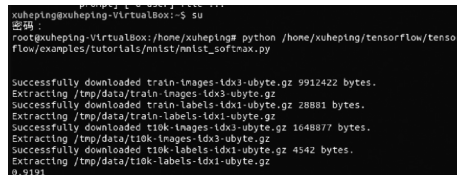


图6 MNIST_softmax.py运行结果

这只是最简单的一个模型,有许多算法模型的正确率可以达到0.997左右。

4 结论

TensorFlow是一个很好的利用机器学习算法的框架,而它的优势在于深度学习系统的构建,虽然在本文中没有涉及,但是从实验仿真中可以看到TensorFlow的模型构建简便,训练速度快。

参考文献

- [1] HARRINGTON P. 机器学习实战[M]. 李锐,李鹏,曲亚东,等,译. 北京:人民邮电出版社,2013.
- [2] TensorFlow官方文档中文版[EB/OL]. (2015-11-18) [2016-11-25] <http://wiki.jikexueyuan.com/project/tensorflow-zh/>.
- [3] TensorFlow官方网站[EB/OL]. [2016-11-25] <https://www.tensorflow.org/>.
- [4] TensorFlow架构[EB/OL]. (2016-06-12) [2016-11-25] <http://blog.csdn.net/stdcoutzyx/article/details/51645396>.
- [5] Google TensorFlow机器学习框架介绍和使用[EB/OL]. (2015-12-15) [2016-11-25] http://blog.csdn.net/sinat_31628525/article/details/50320817.
- [6] 张俊,李鑫. TensorFlow平台下的手写字符识别[J]. 电脑知识及技术,2016,12(16):199-201.

(收稿日期:2016-11-25)

作者简介:

章敏敏(1996-),女,本科,主要研究方向:机器学习理论及应用。

徐和平(1995-),男,本科,主要研究方向:机器学习理论及应用。

王晓洁(1996-),女,本科,主要研究方向:机器学习理论及应用。

(上接第57页)

- [5] Chen Xiaoyun, Jian Cairen. Gene expression data clustering based on graph regularized subspace segmentation[J]. Neuro-computing, 2014, 143(16):44-50.
- [6] 林莉媛,陈晓云,简彩仁. 融入距离信息的最小二乘回归子空间分割[J]. 微型机与应用,2016,35(6):63-65.
- [7] VIDAL R. A tutorial on subspace clustering[J]. IEEE Signal Processing Magazine, 2010, 28(2):52-68.
- [8] Shi Jianbo, MALIK J. Normalized cuts and image segmentation[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2000, 22(8):888-905.

- [9] Cai Deng, He Xiaofei, Wu Xiaoyun, et al. Non-negative matrix factorization on manifold[C]. Proceedings of the 8th IEEE International Conference on Data Mining, Pisa, Italy, 2008:63-72.

(收稿日期:2016-11-18)

作者简介:

简彩仁(1988-),男,硕士,主要研究方向:数据挖掘、模式识别等。

吕书龙(1977-),男,硕士,副教授,主要研究方向:数据挖掘、统计计算、应用统计分析等。

《微型机与应用》2017年第36卷第10期