

Unsupervised Feature Learning and Deep Learning

Edited by 陈超 L^AT_EX
chaochancs@gmail.com

致谢

此教程来自 Stanford 的 Andrew Ng [http://ufldl.stanford.edu/wiki/index.php/UFLDL_Tutorial]。这是通过众多网友的合力翻译而成。在此对这些网友的无私贡献致意！翻译者如下。

神经网络：

孙逊 (sunpaofufoxmail.com)，林锋 (xlfgyeah.net)，
刘鸿鹏飞 (just.darkfoxmail.com)，许利杰 (csxulijie@gmail.com)。

反向传导算法：

王方 (fangkey@gmail.com)，林锋 (xlfg@yeah.net)，
许利杰 (csxulijie@gmail.com)。

梯度检验与高级优化：

袁晓丹 (shadowwalker1991@gmail.com)，王方 (fangkey@gmail.com)，
林锋 (xlfg@yeah.net)，许利杰 (csxulijie@gmail.com)。

自编码算法与稀疏性：

周韬 (ztsailing@gmail.com)，葛燕儒 (yrgehi@gmail.com)，
林锋 (xlfg@yeah.net)，余凯 (kai.yu.cool@gmail.com)。

可视化自编码器训练结果：

王方 (fangkey@gmail.com)，胡伦 (hulun499@gmail.com)，
谢宇 (msforbus@sina.com)，@小琳爱肉肉 (新浪微博账号)，
余凯 (kai.yu.cool@gmail.com)。

稀疏自编码器符号一览表：

邵杰 (jiesha@hotmail.com)，许利杰 (csxulijie@gmail.com)，
余凯 (kai.yu.cool@gmail.com)。

矢量化编程：

彭君睿 (07caleb@gmail.com)，王文中 (wangwenzhong@ymail.com)，
邓亚峰 (dengyafeng@gmail.com)。

逻辑回归的向量化实现样例：

林锋 (xlfg@yeah.net)，谭晓阳 (x.tan@nuaa.edu.cn)，
邓亚峰 (dengyafeng@gmail.com)。

神经网络向量化：

阎志涛 (zhitaoyan@gmail.com)，谭晓阳 (x.tan@nuaa.edu.cn)，
邓亚峰 (dengyafeng@gmail.com)。

主成分分析：

郭亮 (guoliang2248@gmail.com)，张力 (emma.lzhang@gmail.com)，
金峰 (jinfengb@gmail.com)，@破破的桥 (新浪微博)，谭晓阳 (x.tan@nuaa.edu.cn)。

白化：

杨海川 (yanghaichuan@outlook.com)，王文中 (wangwenzhong@ymail.com)，
谭晓阳 (x.tan@nuaa.edu.cn)。

实现主成分分析和白化：

周思远 (visualzhou@gmail.com)，张力 (emma.lzhang@gmail.com)，
谭晓阳 (x.tan@nuaa.edu.cn)。

Softmax回归:

曾俊瑀 (knighterzjy@gmail.com), 王方 (fangkey@gmail.com),
王文中 (wangwenzhong@ymail.com)。

自我学习:

曾俊瑀 (knighterzjy@gmail.com), 王方 (fangkey@gmail.com),
王文中 (wangwenzhong@ymail.com)。

从自我学习到深层网络:

杨耀 (iamyangyao@163.com), 阎志涛 (zhitao.yan@gmail.com),
王文中 (wangwenzhong@ymail.com)。

深度网络概览

郑胤 (yzheng3xg@gmail.com), 谭晓阳 (x.tan@nuaa.edu.cn),
许利杰 (csxulijie@gmail.com)。

栈式自编码算法:

张天雷 (ztl2004@gmail.com), 邓亚峰 (dengyafeng@gmail.com),
许利杰 (csxulijie@gmail.com)。

微调多层自编码算法:

崔巍 (watsoncui@gmail.com), 余凯 (kai.yu.cool@gmail.com),
许利杰 (csxulijie@gmail.com)。

线性解码器:

严晓东 (yan.endless@gmail.com), 姚涛 (yaothinker@gmail.com),
晓风 (xiaofeng.zhb@alibaba-inc.com)。

卷积特征提取:

孔德硕 (tobluestone@gmail.com), 郭亮 (guoliang2248@gmail.com),
晓风 (xiaofeng.zhb@alibaba-inc.com)。

池化:

陈玉栓 (chris_chen_cys@hotmail.com), 刘鸿鹏飞 (just.dark@foxmail.com),
邓亚峰 (dengyafeng@gmail.com), 晓风 (xiaofeng.zhb@alibaba-inc.com)。

再次感谢以上贡献者!

Chapter 1

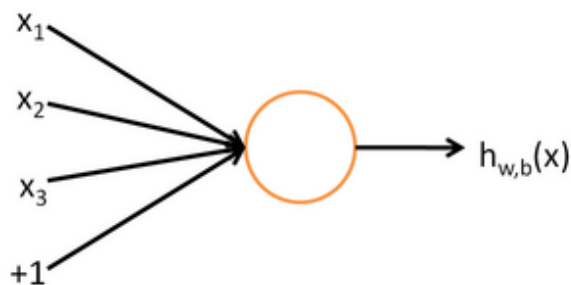
稀疏自编码器

1.1 神经网络

1.1.1 概述

以监督学习为例，假设我们有训练样本集 $(x^{(i)}, y^{(i)})$ ，那么神经网络算法能够提供一种复杂且非线性的假设模型 $h_{W,b}(x)$ ，它具有参数 W, b ，可以以此参数来拟合我们的数据

为了描述神经网络，我们先从最简单的神经网络讲起 这个神经网络仅由一个“神经元”构成，以下即是这个“神经元”的图示：



这个“神经元”是一个以 x_1, x_2, x_3 及截距 $+1$ 为输入值的运算单元，其输出为 $h_{W,b}(x) = f(W^T x) = f(\sum_{i=1}^3 W_i x_i + b)$ ，其中函数 $f: \mathbb{R} \mapsto \mathbb{R}$ 被称为“激活函数”。在本教程中，我们选用sigmoid函数作为激活函数 $f(\cdot)$

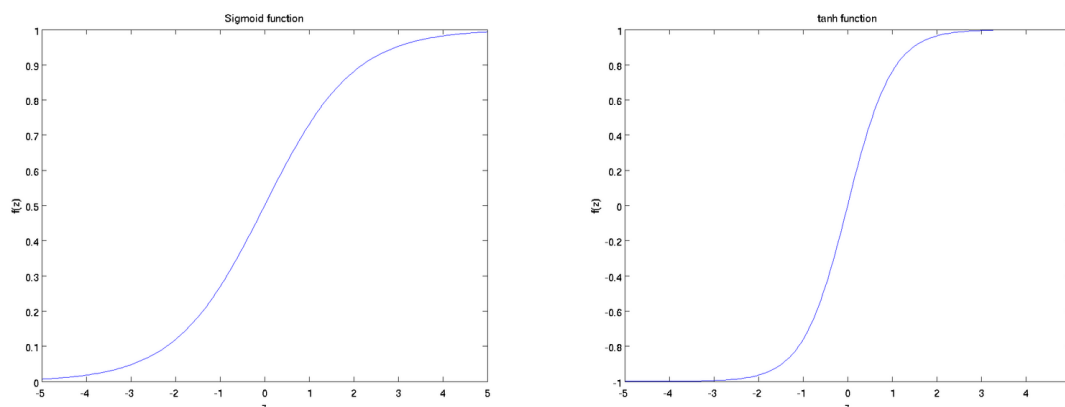
$$f(z) = \frac{1}{1 + \exp(-z)}. \quad (1.1)$$

可以看出，这个单一“神经元”的输入—输出映射关系其实就是

一个逻辑回归（logistic regression）。虽然本系列教程采用sigmoid函数，但你也可以选择双曲正切函数（tanh）：

$$f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}, \quad (1.2)$$

以下分别是sigmoid及tanh的函数图像



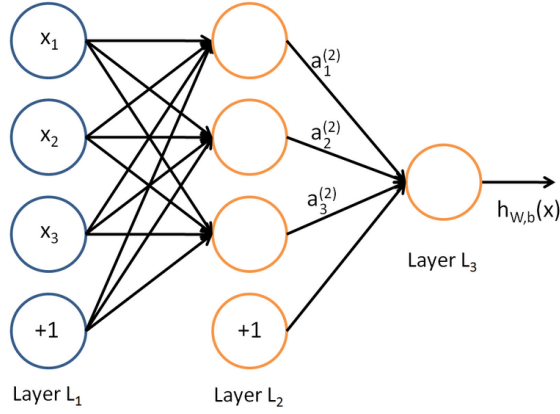
$\tanh(z)$ 函数是sigmoid函数的一种变体，它的取值范围为 $[-1, 1]$ ，而不是sigmoid函数的 $[0, 1]$ 。

注意，与其它地方（包括OpenClassroom公开课以及斯坦福大学CS229课程）不同的是，这里我们不再令 $x_0 = 1$ 。取而代之，我们用单独的参数 b 来表示截距。

最后要说明的是，有一个等式我们以后会经常用到：如果选择 $f(z) = 1/(1 + \exp(-z))$ ，也就是sigmoid函数，那么它的导数就是 $f'(z) = f(z)(1 - f(z))$ （如果选择tanh函数，那它的导数就是 $f'(z) = 1 - (f(z))^2$ ，你可以根据sigmoid（或tanh）函数的定义自行推导这个等式。

1.1.2 神经网络模型

所谓神经网络就是将许多个单一“神经元”联结在一起，这样，一个“神经元”的输出就可以是另一个“神经元”的输入。例如，下图就是一个简单的神经网络：



我们使用圆圈来表示神经网络的输入，标上“+1”的圆圈被称为偏置节点，也就是截距项。神经网络最左边的一层叫做输入层，最右的一层叫做输出层（本例中，输出层只有一个节点）。中间所有节点组成的一层叫做隐藏层，因为我们不能在训练样本集中观测到它们的值。同时可以看到，以上神经网络的例子中有3个输入单元（偏置单元不计在内），3个隐藏单元及一个输出单元。

我们用 n_l 来表示网络的层数，本例中 $n_l = 3$ ，我们将第 l 层记为 L_l ，于是 L_1 是输入层，输出层是 L_{n_l} 。本例神经网络有参数 $(W, b) = (W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)})$ ，其中 $W_{ij}^{(l)}$ （下面的式子中用到）是第 l 层第 j 单元与第 $l+1$ 层第 i 单元之间的联接参数（其实就是连接线上的权重，注意标号顺序）， $b_i^{(l)}$ 是第 $l+1$ 层第 i 单元的偏置项。因此在本例中， $W^{(1)} \in \Re^{3 \times 3}$ ， $W^{(2)} \in \Re^{1 \times 3}$ 。注意，没有其他单元连向偏置单元（即偏置单元没有输入），因为它们总是输出 +1。同时，我们用 s_l 表示第 l 层的节点数（偏置单元不计在内）。

我们用 $a_i^{(l)}$ 表示第 l 层第 i 单元的激活值（输出值）。当 $l=1$ 时， $a_i^{(1)} = x_i$ ，也就是第 i 个输入值（输入值的第 i 个特征）。对于给定参数集合 W, b ，我们的神经网络就可以按照函数 $h_{W,b}(x)$ 来计算输出结果。本例神经网络的计算步骤如下：

$$a_1^{(2)} = f(W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)}) \quad (1.3)$$

$$a_2^{(2)} = f(W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{23}^{(1)}x_3 + b_2^{(1)}) \quad (1.4)$$

$$a_3^{(2)} = f(W_{31}^{(1)}x_1 + W_{32}^{(1)}x_2 + W_{33}^{(1)}x_3 + b_3^{(1)}) \quad (1.5)$$

$$h_{W,b}(x) = a_1^{(3)} = f(W_{11}^{(2)}a_1^{(2)} + W_{12}^{(2)}a_2^{(2)} + W_{13}^{(2)}a_3^{(2)} + b_1^{(2)}) \quad (1.6)$$

我们用 $z_i^{(l)}$ 表示第 l 层第 i 单元输入加权和（包括偏置单元），比如， $z_i^{(2)} = \sum_{j=1}^n W_{ij}^{(1)}x_j + b_i^{(1)}$ ，则 $a_i^{(l)} = f(z_i^{(l)})$ 。这样我们就可以得到一种更简洁的表示法。这里我们将激活函数 $f(\cdot)$ 扩展为用向量（分量的形式）来表示，即 $f([z_1, z_2, z_3]) = [f(z_1), f(z_2), f(z_3)]$ ，那么，上面的等式可以更简洁地表示为：

$$z^{(2)} = W^{(1)}x + b^{(1)} \quad (1.7)$$

$$a^{(2)} = f(z^{(2)}) \quad (1.8)$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)} \quad (1.9)$$

$$h_{W,b}(x) = a^{(3)} = f(z^{(3)}) \quad (1.10)$$

我们将上面的计算步骤叫作前向传播。回想一下，之前我们用 $a^{(1)} = x$ 表示输入层的激活值，那么给定第 l 层的激活值 $a^{(l)}$ 后，第 $l+1$ 层的激活值 $a^{(l+1)}$ 就可以按照下面步骤计算得到：

$$z^{(l+1)} = W^{(l)}a^{(l)} + b^{(l)} \quad (1.11)$$

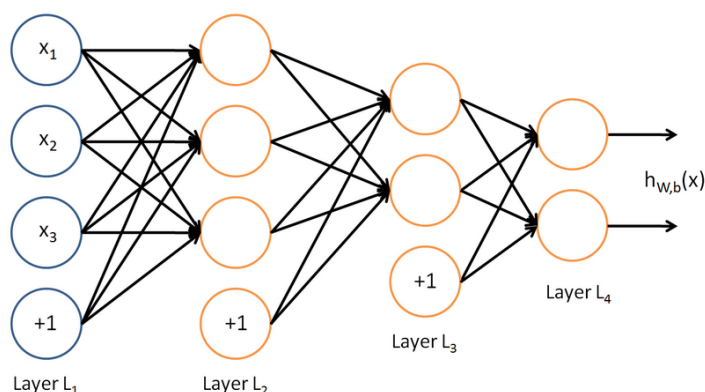
$$a^{(l+1)} = f(z^{(l+1)}) \quad (1.12)$$

将参数矩阵化，使用矩阵一向量运算方式，我们就可以利用线性代数的优势对神经网络进行快速求解。

目前为止，我们讨论了一种神经网络，我们也可以构建另一种结

构的神经网络（这里结构指的是神经元之间的联接模式），也就是包含多个隐藏层的神经网络。最常见的一个例子是 n_l 层的神经网络，第 1 层是输入层，第 n_l 层是输出层，中间的每个层 l 与层 $l+1$ 紧密相联。这种模式下，要计算神经网络的输出结果，我们可以按照之前描述的等式，按部就班，进行前向传播，逐一计算第 L_2 层的所有激活值，然后是第 L_3 层的激活值，以此类推，直到第 L_{n_l} 层。这是一个前馈神经网络的例子，因为这种联接图没有闭环或回路。

神经网络也可以有多个输出单元。比如，下面的神经网络有两层隐藏层： L_2 及 L_3 ，输出层 L_4 有两个输出单元。



要求解这样的神经网络，需要样本集 $(x^{(i)}, y^{(i)})$ ，其中 $y^{(i)} \in \mathbb{R}^2$ 。如果你想预测的输出是多个的，那这种神经网络很适用。（比如，在医疗诊断应用中，患者的体征指标就可以作为向量的输入值，而不同的输出值 y_i 可以表示不同的疾病存在与否。）

1.2 反向传导算法

假设我们有一个固定样本集 $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ ，它包含 m 个样例。我们可以用批量梯度下降法来求解神经网络。具体来讲，对于单个样例 (x, y) ，其代价函数为：

$$J(W, b; x, y) = \frac{1}{2} \|h_{W,b}(x) - y\|^2. \quad (1.13)$$

这是一个（二分之一的）方差代价函数。给定一个包含 m 个样例的数据集，我们可以定义整体代价函数为：

$$\begin{aligned} J(W, b) &= \left[\frac{1}{m} \sum_{i=1}^m J(W, b; x^{(i)}, y^{(i)}) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} \left(W_{ji}^{(l)} \right)^2 \quad (1.14) \\ &= \left[\frac{1}{m} \sum_{i=1}^m \left(\frac{1}{2} \|h_{W,b}(x^{(i)}) - y^{(i)}\|^2 \right) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} \left(W_{ji}^{(l)} \right)^2 \quad (1.15) \end{aligned}$$

以上公式中的第一项 $J(W, b)$ 是一个均方差项。第二项是一个规则化项（也叫权重衰减项），其目的是减小权重的幅度，防止过度拟合。

[注：通常权重衰减的计算并不使用偏置项 $b_i^{(l)}$ ，比如我们在 $J(W, b)$ 的定义中就没有使用。一般来说，将偏置项包含在权重衰减项中只会对最终的神经网络产生很小的影响。如果你在斯坦福选修过CS229（机器学习）课程，或者在YouTube上看过课程视频，你会发现这个权重衰减实际上是课上提到的贝叶斯规则化方法的变种。在贝叶斯规则化方法中，我们将高斯先验概率引入到参数中计算MAP（极大后验）估计（而不是极大似然估计）。]

权重衰减参数 λ 用于控制公式中两项的相对重要性。在此重申一下这两个复杂函数的含义： $J(W, b; x, y)$ 是针对单个样例计算得到的方差代价函数； $J(W, b)$ 是整体样本代价函数，它包含权重衰减项。

以上的代价函数经常被用于分类和回归问题。在分类问题中，我们用 $y = 0$ 或 1 ，来代表两种类型的标签（回想一下，这是因为 sigmoid 激活函数的值域为 $[0, 1]$ ；如果我们使用双曲正切型激活函数，那么应该选用 -1 和 $+1$ 作为标签）。对于回归问题，我们首

先要变换输出值域（译者注：也就是 y ），以保证其范围为 $[0, 1]$ （同样地，如果我们使用双曲正切型激活函数，要使输出值域为 $[-1, 1]$ ）。

我们的目标是针对参数 W 和 b 来求其函数 $J(W, b)$ 的最小值。为了求解神经网络，我们需要将每一个参数 $W_{ij}^{(l)}$ 和 $b_i^{(l)}$ 初始化为一个很小的、接近零的随机值（比如说，使用正态分布 $Normal(0, \epsilon^2)$ 生成的随机值，其中 ϵ 设置为 0.01，之后对目标函数使用诸如批量梯度下降法的最优化算法。因为 $J(W, b)$ 是一个非凸函数，梯度下降法很可能会收敛到局部最优解；但是在实际应用中，梯度下降法通常能得到令人满意的结果。最后，需要再次强调的是，要将参数进行随机初始化，而不是全部置为 0。如果所有参数都用相同的值作为初始值，那么所有隐藏层单元最终会得到与输入值有关的、相同的函数（也就是说，对于所有 i ， $W_{ij}^{(1)}$ 都会取相同的值，那么对于任何输入 x 都会有： $a_1^{(2)} = a_2^{(2)} = a_3^{(2)} = \dots \Psi$ 。随机初始化的目的是使对称失效。

梯度下降法中每一次迭代都按照如下公式对参数 W 和 b 进行更新：

$$W_{ij}^{(l)} = W_{ij}^{(l)} - \alpha \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b) \quad (1.16)$$

$$b_i^{(l)} = b_i^{(l)} - \alpha \frac{\partial}{\partial b_i^{(l)}} J(W, b) \quad (1.17)$$

其中 α 是学习速率。其中关键步骤是计算偏导数。我们现在来讲一下反向传播算法，它是计算偏导数的一种有效方法。

我们首先来讲一下如何使用反向传播算法来计算 $\frac{\partial}{\partial W_{ij}^{(l)}} J(W, b; x, y)$ 和 $\frac{\partial}{\partial b_i^{(l)}} J(W, b; x, y)$ ，这两项是单个样例 (x, y) 的代价函数 $J(W, b; x, y)$ 的偏导数。一旦我们求出该偏导数，就可以推导出整体代价函数

$J(W, b)$ 的偏导数:

$$\frac{\partial}{\partial W_{ij}^{(l)}} J(W, b) = \left[\frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b; x^{(i)}, y^{(i)}) \right] + \lambda W_{ij}^{(l)} \quad (1.18)$$

$$\frac{\partial}{\partial b_i^{(l)}} J(W, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial b_i^{(l)}} J(W, b; x^{(i)}, y^{(i)}) \quad (1.19)$$

以上两行公式稍有不同，第一行比第二行多出一项，是因为权重衰减是作用于 W 而不是 b 。

反向传播算法的思路如下：给定一个样例 (x, y) ，我们首先进行“前向传导”运算，计算出网络中所有的激活值，包括 $h_{W,b}(x)$ 的输出值。之后，针对第 l 层的每一个节点 i ，我们计算出其“残差” $\delta_i^{(l)}$ ，该残差表明了该节点对最终输出值的残差产生了多少影响。对于最终的输出节点，我们可以直接算出网络产生的激活值与实际值之间的差距，我们将这个差距定义为 $\delta_i^{(n_l)}$ （第 n_l 层表示输出层）。对于隐藏单元我们如何处理呢？我们将基于节点（译者注：第 $l+1$ 层节点）残差的加权平均值计算 $\delta_i^{(l)}$ ，这些节点以 $a_i^{(l)}$ 作为输入。下面将给出反向传导算法的细节：

进行前馈传导计算，利用前向传导公式，得到 L_2, L_3, \dots 直到输出层 L_{n_l} 的激活值。

对于第 n_l 层（输出层）的每个输出单元 i ，我们根据以下公式计算残差：

$$\delta_i^{(n_l)} = \frac{\partial}{\partial z_i^{(n_l)}} \frac{1}{2} \|y - h_{W,b}(x)\|^2 = -(y_i - a_i^{(n_l)}) \cdot f'(z_i^{(n_l)}) \quad (1.20)$$

[译者注:

$$\delta_i^{(n_l)} = \frac{\partial}{\partial z_i^{n_l}} J(W, b; x, y) = \frac{\partial}{\partial z_i^{n_l}} \frac{1}{2} \|y - h_{W,b}(x)\|^2 \quad (1.21)$$

$$= \frac{\partial}{\partial z_i^{n_l}} \frac{1}{2} \sum_{j=1}^{S_{n_l}} (y_j - a_j^{(n_l)})^2 = \frac{\partial}{\partial z_i^{n_l}} \frac{1}{2} \sum_{j=1}^{S_{n_l}} (y_j - f(z_j^{(n_l)}))^2 \quad (1.22)$$

$$= -(y_i - f(z_i^{(n_l)})) \cdot f'(z_i^{(n_l)}) = -(y_i - a_i^{(n_l)}) \cdot f'(z_i^{(n_l)}) \quad (1.23)$$

]

对 $l = n_l - 1, n_l - 2, n_l - 3, \dots, 2$ 的各个层, 第 l 层的第 i 个节点的残差计算方法如下:

$$\delta_i^{(l)} = \left(\sum_{j=1}^{s_{l+1}} W_{ji}^{(l)} \delta_j^{(l+1)} \right) f'(z_i^{(l)})$$

[译者注:

$$\delta_i^{(n_l-1)} = \frac{\partial}{\partial z_i^{n_l-1}} J(W, b; x, y) = \frac{\partial}{\partial z_i^{n_l-1}} \frac{1}{2} \|y - h_{W,b}(x)\|^2 = \frac{\partial}{\partial z_i^{n_l-1}} \frac{1}{2} \sum_{j=1}^{S_{n_l}} (y_j - a_j^{(n_l)})^2 \quad (1.24)$$

$$= \frac{1}{2} \sum_{j=1}^{S_{n_l}} \frac{\partial}{\partial z_i^{n_l-1}} (y_j - a_j^{(n_l)})^2 = \frac{1}{2} \sum_{j=1}^{S_{n_l}} \frac{\partial}{\partial z_i^{n_l-1}} (y_j - f(z_j^{(n_l)}))^2 \quad (1.25)$$

$$= \sum_{j=1}^{S_{n_l}} -(y_j - f(z_j^{(n_l)})) \cdot \frac{\partial}{\partial z_i^{(n_l-1)}} f(z_j^{(n_l)}) \quad (1.26)$$

$$= \sum_{j=1}^{S_{n_l}} \delta_j^{(n_l)} \cdot \frac{\partial z_j^{(n_l)}}{\partial z_i^{n_l-1}} = \sum_{j=1}^{S_{n_l}} \left(\delta_j^{(n_l)} \cdot \frac{\partial}{\partial z_i^{n_l-1}} \sum_{k=1}^{S_{n_l-1}} f(z_k^{n_l-1}) \cdot W_{jk}^{n_l-1} \right) \quad (1.27)$$

$$= \sum_{j=1}^{S_{n_l}} \delta_j^{(n_l)} \cdot W_{ji}^{n_l-1} \cdot f'(z_i^{n_l-1}) = \left(\sum_{j=1}^{S_{n_l}} W_{ji}^{n_l-1} \delta_j^{(n_l)} \right) f'(z_i^{n_l-1}) \quad (1.28)$$

将上式中的 $n_l - 1$ 与 n_l 的关系替换为 l 与 $l + 1$ 的关系，就可以得到：

$\delta_i^{(l)} = \left(\sum_{j=1}^{S_{l+1}} W_{ji}^{(l)} \delta_j^{(l+1)} \right) f'(z_i^{(l)})$ 以上逐次从后向前求导的过程即为“反向传导”的本意所在。]

计算我们需要的偏导数，计算方法如下：

$$\frac{\partial}{\partial W_{ij}^{(l)}} J(W, b; x, y) = a_j^{(l)} \delta_i^{(l+1)} \quad (1.29)$$

$$\frac{\partial}{\partial b_i^{(l)}} J(W, b; x, y) = \delta_i^{(l+1)}. \quad (1.30)$$

最后，我们用矩阵-向量表示法重写以上算法。我们使用“ \bullet ”表示

向量乘积运算符（在Matlab或Octave里用“.”表示，也称作阿达马乘积）。若 $a = b \bullet c$ ，则 $a_i = b_i c_i$ 。在上一个教程中我们扩展了 $f(\cdot)$ 的定义，使其包含向量运算，这里我们也对偏导数 $f'(\cdot)$ 也做了同样的处理（于是又有 $f'([z_1, z_2, z_3]) = [f'(z_1), f'(z_2), f'(z_3)]\Psi$ 。

那么，反向传播算法可表示为以下几个步骤：进行前馈传导计算，利用前向传导公式，得到 L_2, L_3, \dots 直到输出层 L_{n_l} 的激活值。对输出层（第 n_l 层），计算：

$$\delta^{(n_l)} = -(y - a^{(n_l)}) \bullet f'(z^{(n_l)}) \quad (1.31)$$

对于 $l = n_l - 1, n_l - 2, n_l - 3, \dots, 2$ 的各层，计算：

$$\delta^{(l)} = \left((W^{(l)})^T \delta^{(l+1)} \right) \bullet f'(z^{(l)}) \quad (1.32)$$

计算最终需要的偏导数值：

$$\nabla_{W^{(l)}} J(W, b; x, y) = \delta^{(l+1)} (a^{(l)})^T, \quad (1.33)$$

$$\nabla_{b^{(l)}} J(W, b; x, y) = \delta^{(l+1)}. \quad (1.34)$$

实现中应注意：在以上的第2步和第3步中，我们需要为每一个 i 值计算其 $f'(z_i^{(l)})$ 。假设 $f(z)$ 是sigmoid函数，并且我们已经在前向传导运算中得到了 $a_i^{(l)}$ 。那么，使用我们早先推导出的 $f'(z)$ 表达式，就可以计算得到 $f'(z_i^{(l)}) = a_i^{(l)}(1 - a_i^{(l)})$ 。

最后，我们将对梯度下降算法做个全面总结。在下面的伪代码中， $\Delta W^{(l)}$ 是一个与矩阵 $W^{(l)}$ 维度相同的矩阵， $\Delta b^{(l)}$ 是一个与 $b^{(l)}$ 维度相同的向量。注意这里“ $\Delta W^{(l)}$ ”是一个矩阵，而不是“ Δ 与 $W^{(l)}$ 相乘”。下面，我们实现批量梯度下降法中的一次迭代：

对于所有 l ，令 $\Delta W^{(l)} := 0, \Delta b^{(l)} := 0$ （设置为全零矩阵或全零向量）对于 $i = 1$ 到 m ，使用反向传播算法计算 $\nabla_{W^{(l)}} J(W, b; x, y)$ 和 $\nabla_{b^{(l)}} J(W, b; x, y)$ 。计算 $\Delta W^{(l)} := \Delta W^{(l)} + \nabla_{W^{(l)}} J(W, b; x, y)$ 。计算

$\Delta b^{(l)} := \Delta b^{(l)} + \nabla_{b^{(l)}} J(W, b; x, y)$ 。更新权重参数：

$$W^{(l)} = W^{(l)} - \alpha \left[\left(\frac{1}{m} \Delta W^{(l)} \right) + \lambda W^{(l)} \right] \quad (1.35)$$

$$b^{(l)} = b^{(l)} - \alpha \left[\frac{1}{m} \Delta b^{(l)} \right] \quad (1.36)$$

现在，我们可以重复梯度下降法的迭代步骤来减小代价函数 $J(W, b)$ 的值，进而求解我们的神经网络。

1.3 梯度检验与高级优化

众所周知，反向传播算法很难调试得到正确结果，尤其是当实现程序存在很多难于发现的bug时。举例来说，索引的缺位错误（off-by-one error）会导致只有部分层的权重得到训练，再比如忘记计算偏置项。这些错误会使你得到一个看似十分合理的结果（但实际上比正确代码的结果要差）。因此，但从计算结果上来看，我们很难发现代码中有什么东西遗漏了。本节中，我们将介绍一种对求导结果进行数值检验的方法，该方法可以验证求导代码是否正确。另外，使用本节所述求导检验方法，可以帮助你提升写正确代码的信心。

缺位错误（Off-by-one error）举例说明：比如 `for` 循环中循环 m 次，正确应该是 `for(i = 1; i <= m; i++)`，但有时程序员疏忽，会写成 `for(i = 1; i < m; i++)`，这就是缺位错误。

假设我们想要最小化以 θ 为自变量的目标函数 $J(\theta)$ 。假设 $J : \Re \mapsto \Re$ ，则 $\theta \in \Re$ 。在一维的情况下，一次迭代的梯度下降公式是

$$\theta := \theta - \alpha \frac{d}{d\theta} J(\theta). \quad (1.37)$$

再假设我们已经用代码实现了计算 $\frac{d}{d\theta} J(\theta)$ 的函数 $g(\theta)$ ，接着我们使用 $\theta := \theta - \alpha g(\theta)$ 来实现梯度下降算法。那么我们如何检验 g 的

实现是否正确呢？回忆导数的数学定义：

$$\frac{d}{d\theta} J(\theta) = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}. \quad (1.38)$$

那么对于任意 θ 值，我们都可以对等式左边的导数用：

$$\frac{J(\theta + \text{EPSILON}) - J(\theta - \text{EPSILON})}{2 \times \text{EPSILON}} \quad (1.39)$$

来近似。

实际应用中，我们常将 *EPSILON* 设为一个很小的常量，比如在 10^{-4} 数量级（虽然 *EPSILON* 的取值范围可以很大，但是我们不会将它设得太小，比如 10^{-20} ，因为那将导致数值舍入误差。）

给定一个被认为能计算 $\frac{d}{d\theta} J(\theta)$ 的函数 $g(\theta)$ ，我们可以用下面的数值检验公式

$$g(\theta) \approx \frac{J(\theta + \text{EPSILON}) - J(\theta - \text{EPSILON})}{2 \times \text{EPSILON}}. \quad (1.40)$$

计算两端是否一样来检验函数是否正确。

上式两端值的接近程度取决于 J 的具体形式。但是在假定 $\text{EPSILON} = 10^{-4}$ 的情况下，你通常会发现上式左右两端至少有4位有效数字是一样的（通常会更多）。

现在，考虑 $\theta \in \mathbb{R}^n$ 是一个向量而非一个实数（那么就有 n 个参数要学习得到），并且 $J: \mathbb{R}^n \mapsto \mathbb{R}$ 。在神经网络的例子里我们使用 $J(W, b)$ ，可以想象为把参数 W, b 组合扩展成一个长向量 θ 。现在我们将求导检验方法推广到一般化，即 θ 是一个向量的情况。假设我们有一个用于计算 $\frac{\partial}{\partial \theta_i} J(\theta)$ 的函数 $g_i(\theta)$ ；我们想要检验 g_i 是否输

出正确的求导结果。我们定义 $\theta^{(i+)} = \theta + \text{EPSILON} \times \vec{e}_i$ ，其中

$$\vec{e}_i = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad (1.41)$$

是第 i 个基向量（维度和 θ 相同，在第 i 行是“1”而其他行是“0”）。所以， $\theta^{(i+)}$ 和 θ 几乎相同，除了第 i 行元素增加了 EPSILON 。类似地， $\theta^{(i-)} = \theta - \text{EPSILON} \times \vec{e}_i$ 得到的第 i 行减小了 EPSILON 。然后我们可以对每个 i 检查下式是否成立，进而验证 $g_i(\theta)$ 的正确性：

$$g_i(\theta) \approx \frac{J(\theta^{(i+)}) - J(\theta^{(i-)})}{2 \times \text{EPSILON}}. \quad (1.42)$$

当用反射传播算法求解神经网络时，正确算法实现会得到：

$$\nabla_{W^{(l)}} J(W, b) = \left(\frac{1}{m} \Delta W^{(l)} \right) + \lambda W^{(l)} \quad (1.43)$$

$$\nabla_{b^{(l)}} J(W, b) = \frac{1}{m} \Delta b^{(l)}. \quad (1.44)$$

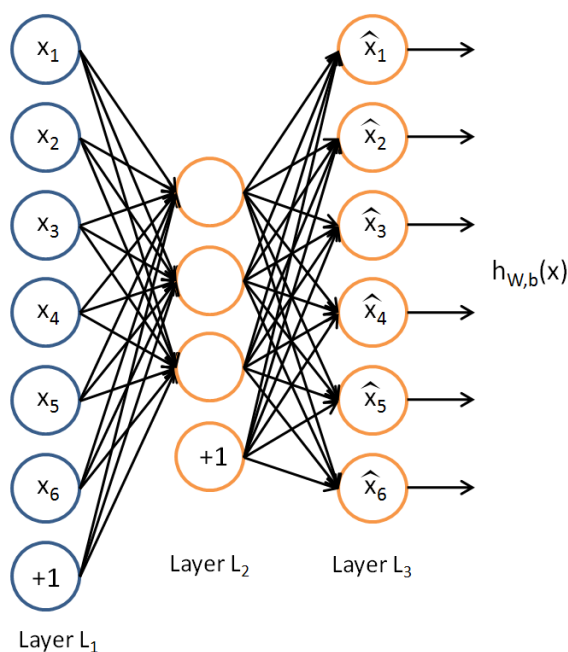
以上结果与反向传播算法中的最后一段伪代码一致，都是计算梯度下降。为了验证梯度下降代码的正确性，使用上述数值检验方法计算 $J(W, b)$ 的导数，然后验证 $\left(\frac{1}{m} \Delta W^{(l)} \right) + \lambda W$ 与 $\frac{1}{m} \Delta b^{(l)}$ 是否能够给出正确的求导结果。

迄今为止，我们的讨论都集中在使用梯度下降法来最小化 $J(\theta)$ 。如果你已经实现了一个计算 $J(\theta)$ 和 $\nabla_{\theta} J(\theta)$ 的函数，那么其实还有更精妙的算法来最小化 $J(\theta)$ 。举例来说，可以想象这样一个算法：它使用梯度下降，并能够自动调整学习速率 α ，以得到合适的步长值，最终使 θ 能够快速收敛到一个局部最优解。还有更妙的算法：

比如可以寻找一个Hessian矩阵的近似，得到最佳步长值，使用该步长值能够更快地收敛到局部最优（和牛顿法类似）。此类算法的详细讨论已超出了这份讲义的范围，但是L-BFGS算法我们以后会有论述（另一个例子是共轭梯度算法）。你将在编程练习里使用这些算法中的一个。使用这些高级优化算法时，你需要提供关键的函数：即对于任一个 θ ，需要你计算出 $J(\theta)$ 和 $\nabla_{\theta}J(\theta)$ 。之后，这些优化算法会自动调整学习速率/步长值 α 的大小（并计算Hessian近似矩阵等等）来自动寻找 $J(\theta)$ 最小化时 θ 的值。诸如L-BFGS和共轭梯度算法通常比梯度下降法快很多。

1.4 自编码算法与稀疏性

目前为止，我们已经讨论了神经网络在有监督学习中的应用。在有监督学习中，训练样本是有类别标签的。现在假设我们只有一个没有带类别标签的训练样本集合 $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots\}$ ，其中 $x^{(i)} \in \mathbb{R}^n$ 。自编码神经网络是一种无监督学习算法，它使用了反向传播算法，并让目标值等于输入值，比如 $y^{(i)} = x^{(i)}$ 。下图是一个自编码神经网络的示例。



自编码神经网络尝试学习一个 $h_{W,b}(x) \approx x$ 的函数。换句话说，它尝试逼近一个恒等函数，从而使得输出 \hat{x} 接近于输入 x 。恒等函数虽然看上去不太有学习的意义，但是当我们为自编码神经网络加入某些限制，比如限定隐藏神经元的数量，我们就可以从输入数据中发现一些有趣的结构。举例来说，假设某个自编码神经网络的输入 x 是一张 10×10 图像（共100个像素）的像素灰度值，于是 $n = 100$ ，其隐藏层 L_2 中有50个隐藏神经元。注意，输出也是100维的 $y \in \mathbb{R}^{100}$ 。由于只有50个隐藏神经元，我们迫使自编码神经网络去学习输入数据的压缩表示，也就是说，它必须从50维的隐藏神经元激活度向量 $a^{(2)} \in \mathbb{R}^{50}$ 中重构出100维的像素灰度值输入 x 。如果网络的输入数据是完全随机的，比如每一个输入 x_i 都是一个跟其它特征完全无关的独立同分布高斯随机变量，那么这一压缩表示将会非常难学习。但是如果输入数据中隐含着一些特定的结构，比如某些输入特征是彼此相关的，那么这一算法就可以发现输入数据中的这些相关性。事实上，这一简单的自编码神经网络通常可以学习出一个跟主元分析（PCA）结果非常相似的输入数据的低维表示。

我们刚才的论述是基于隐藏神经元数量较小的假设。但是即使隐藏神经元的数量较大（可能比输入像素的个数还要多），我们仍然通过给自编码神经网络施加一些其他的限制条件来发现输入数据中的结构。具体来说，如果我们给隐藏神经元加入稀疏性限制，那么自编码神经网络即使在隐藏神经元数量较多的情况下仍然可以发现输入数据中一些有趣的结构。稀疏性可以被简单地解释如下。如果当神经元的输出接近于1的时候我们认为它被激活，而输出接近于0的时候认为它被抑制，那么使得神经元大部分的时间都是被抑制的限制则被称作稀疏性限制。这里我们假设的神经元的激活函数是sigmoid函数。如果你使用tanh作为激活函数的话，当神经元输出

为-1的时候，我们认为神经元是被抑制的。注意到 $a_j^{(2)}$ 表示隐藏神经元 j 的激活度，但是这一表示方法中并未明确指出哪一个输入 x 带来了这一激活度。所以我们将使用 $a_j^{(2)}(x)$ 来表示在给定输入为 x 情况下，自编码神经网络隐藏神经元 j 的激活度。进一步，让

$$\hat{\rho}_j = \frac{1}{m} \sum_{i=1}^m [a_j^{(2)}(x^{(i)})] \quad (1.45)$$

表示隐藏神经元 j 的平均活跃度（在训练集上取平均）。我们可以近似的加入一条限制

$$\hat{\rho}_j = \rho, \quad (1.46)$$

其中， ρ 是稀疏性参数，通常是一个接近于0的较小的值（比如 $\rho = 0.05$ ）。换句话说，我们想要让隐藏神经元 j 的平均活跃度接近0.05。为了满足这一条件，隐藏神经元的活跃度必须接近于0。为了实现这一限制，我们将会在我们的优化目标函数中加入一个额外的惩罚因子，而这一惩罚因子将惩罚那些 $\hat{\rho}_j$ 和 ρ 有显著不同的情况从而使得隐藏神经元的平均活跃度保持在较小范围内。惩罚因子的具体形式有很多种合理的选择，我们将会选择以下这一种：

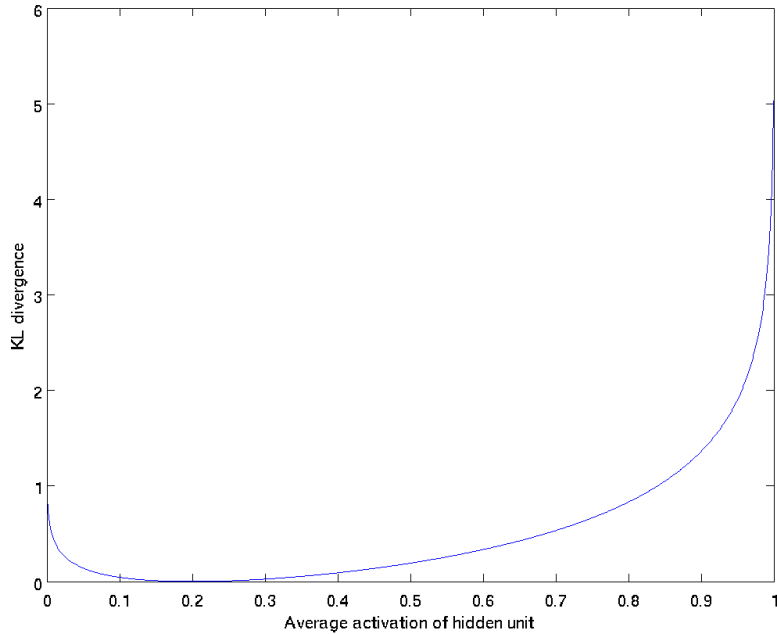
$$\sum_{j=1}^{s_2} \rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_j}. \quad (1.47)$$

这里， s_2 是隐藏层中隐藏神经元的数量，而索引 j 依次代表隐藏层中的每一个神经元。如果你对相对熵（KL divergence）比较熟悉，这一惩罚因子实际上是基于它的。于是惩罚因子也可以被表示为

$$\sum_{j=1}^{s_2} \text{KL}(\rho || \hat{\rho}_j), \quad (1.48)$$

其中 $\text{KL}(\rho || \hat{\rho}_j) = \rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_j}$ 是一个以 ρ 为均值和一个以 $\hat{\rho}_j$ 为均值的两个伯努利随机变量之间的相对熵。相对熵是一种标准的用来测量两个分布之间差异的方法。（如果你没有见过相对熵，不用担心，所有你需要知道的内容都会被包含在这份笔记之中。）

这一惩罚因子有如下性质，当 $\hat{\rho}_j = \rho$ 时 $\text{KL}(\rho||\hat{\rho}_j) = 0$ ，并且随着 $\hat{\rho}_j$ 与 ρ 之间的差异增大而单调递增。举例来说，在下图中，我们设定 $\rho = 0.2$ 并且画出了相对熵值 $\text{KL}(\rho||\hat{\rho}_j)$ 随着 $\hat{\rho}_j$ 变化的变化。



我们可以看出，相对熵在 $\hat{\rho}_j = \rho$ 时达到它的最小值0，而当 $\hat{\rho}_j$ 靠近0或者1的时候，相对熵则变得非常大（其实是趋向于 ∞ ）。所以，最小化这一惩罚因子具有使得 $\hat{\rho}_j$ 靠近 ρ 的效果。现在，我们的总体代价函数可以表示为

$$J_{\text{sparse}}(W, b) = J(W, b) + \beta \sum_{j=1}^{s_2} \text{KL}(\rho||\hat{\rho}_j), \quad (1.49)$$

其中 $J(W, b)$ 如之前所定义，而 β 控制稀疏性惩罚因子的权重。 $\hat{\rho}_j$ 项则也（间接地）取决于 W, b ，因为它是隐藏神经元 j 的平均激活度，而隐藏层神经元的激活度取决于 W, b 。

为了对相对熵进行导数计算，我们可以使用一个易于实现的技巧，这只需要在你的程序中稍作改动即可。具体来说，前面在后向

传播算法中计算第二层（ $l = 2$ ）更新的时候我们已经计算了

$$\delta_i^{(2)} = \left(\sum_{j=1}^{s_2} W_{ji}^{(2)} \delta_j^{(3)} \right) f'(z_i^{(2)}), \quad (1.50)$$

现在我们将其换成

$$\delta_i^{(2)} = \left(\left(\sum_{j=1}^{s_2} W_{ji}^{(2)} \delta_j^{(3)} \right) + \beta \left(-\frac{\rho}{\hat{\rho}_i} + \frac{1 - \rho}{1 - \hat{\rho}_i} \right) \right) f'(z_i^{(2)}). \quad (1.51)$$

就可以了。

有一个需要注意的地方就是我们需要知道 $\hat{\rho}_i$ 来计算这一项更新。所以在计算任何神经元的后向传播之前，你需要对所有的训练样本计算一遍前向传播，从而获取平均激活度。如果你的训练样本可以小到被整个存到内存之中（对于编程作业来说，通常如此），你可以方便地在你所有的样本上计算前向传播并将得到的激活度存入内存并且计算平均激活度。然后你就可以使用事先计算好的激活度来对所有的训练样本进行后向传播的计算。如果你的数据量太大，无法全部存入内存，你就可以扫过你的训练样本并计算一次前向传播，然后将获得的结果累积起来并计算平均激活度 $\hat{\rho}_i$ （当某一个前向传播的结果中的激活度 $a_i^{(2)}$ 被用于计算平均激活度 $\hat{\rho}_i$ 之后就可以将此结果删除）。然后当你完成平均激活度 $\hat{\rho}_i$ 的计算之后，你需要重新对每一个训练样本做一次前向传播从而可以对其进行后向传播的计算。对于后一种情况，你对每一个训练样本需要计算两次前向传播，所以在计算上的效率会稍低一些。证明上面算法能达到梯度下降效果的完整推导过程不再本教程的范围之内。不过如果你想要使用经过以上修改的后向传播来实现自编码神经网络，那么你就会对目标函数 $J_{\text{sparse}}(W, b)$ 做梯度下降。使用梯度验证方法，你可以自己来验证梯度下降算法是否正确。

1.5 可视化自编码器训练结果

训练完（稀疏）自编码器，我们还想把这自编码器学到的函数可视化出来，好弄明白它到底学到了什么。我们以在 10×10 图像（即 $n=100$ ）上训练自编码器为例。在该自编码器中，每个隐藏单元 i 对如下关于输入的函数进行计算：

$$a_i^{(2)} = f \left(\sum_{j=1}^{100} W_{ij}^{(1)} x_j + b_i^{(1)} \right). \quad (1.52)$$

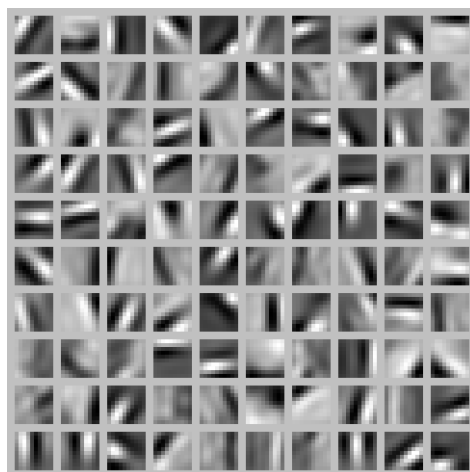
我们将要可视化的函数，就是上面这个以2D图像为输入、并由隐藏单元 i 计算出来的函数。它是依赖于参数 $W_{ij}^{(1)}$ 的（暂时忽略偏置项 b_i ）。需要注意的是， $a_i^{(2)}$ 可看作输入 x 的非线性特征。不过还有个问题：什么样的输入图像 x 可让 $a_i^{(2)}$ 得到最大程度的激励？（通俗一点说，隐藏单元 i 要找个什么样的特征？）。这里我们必须给 x 加约束，否则会得到平凡解。若假设输入有范数约束 $\|x\|^2 = \sum_{i=1}^{100} x_i^2 \leq 1$ ，则可证（请读者自行推导）令隐藏单元 i 得到最大激励的输入应由下面公式计算的像素 x_j 给出（共需计算100个像素， $j=1, \dots, 100$ ）：

$$x_j = \frac{W_{ij}^{(1)}}{\sqrt{\sum_{j=1}^{100} (W_{ij}^{(1)})^2}}. \quad (1.53)$$

当我们用上式算出各像素的值、把它们组成一幅图像、并将图像呈现在我们面前之时，隐藏单元 i 所追寻特征的真正含义也渐渐明朗起来。假如我们训练的自编码器有100个隐藏单元，可视化结果就会包含100幅这样的图像——每个隐藏单元都对应一幅图像。审视这100幅图像，我们可以试着体会这些隐藏单元学出来的整体效果是什么样的。

当我们将稀疏自编码器（100个隐藏单元，在 10×10 像素的输入上

训练) 进行上述可视化处理之后，结果如下所示：



上图的每个小方块都给出了一个（带有有界范数的）输入图像 x ，它可使这100个隐藏单元中的某一个获得最大激励。我们可以看到，不同的隐藏单元学会了在图像的不同位置和方向进行边缘检测。显而易见，这些特征对物体识别等计算机视觉任务是十分有用的。若将其用于其他输入域（如音频），该算法也可学到对这些输入域有用的表示或特征。

1.6 稀疏自编码器符号一览表

下面是我们在推导sparse autoencoder时使用的符号一览表：

符号	含义
x	训练样本的输入特征, $x \in \mathbb{R}^n$.
y	输出值/目标值. 这里 y 可以是向量. 在autoencoder中, $y = x$.
$(x^{(i)}, y^{(i)})$	第 i 个训练样本
$h_{W,b}(x)$	输入为 x 时的假设输出, 其中包含参数 W, b . 该输出应当与目标值 y 具有相同的维数.
$W_{ij}^{(l)}$	连接第 l 层 j 单元和第 $l+1$ 层 i 单元的参数.
$b_i^{(l)}$	第 $l+1$ 层 i 单元的偏置项. 也可以看作是连接第 l 层偏置单元和第 $l+1$ 层 i 单元的参数.
θ	参数向量. 可以认为该向量是通过将参数 W, b 组合展开为一个长的列向量而得到.
$a_i^{(l)}$	网络中第 l 层 i 单元的激活 (输出) 值. 另外, 由于 L_1 层是输入层, 所以 $a_i^{(1)} = x_i$.
$f(\cdot)$	激活函数. 本文中我们使用 $f(z) = \tanh(z)$.
$z_i^{(l)}$	第 l 层 i 单元所有输入的加权和. 因此有 $a_i^{(l)} = f(z_i^{(l)})$.
α	学习率
s_l	第 l 层的单元数目 (不包含偏置单元).
n_l	网络中的层数. 通常 L_1 层是输入层, L_{n_l} 层是输出层.
λ	权重衰减系数.
\hat{x}	对于一个autoencoder, 该符号表示其输出值; 亦即输入值 x 的重构值. 与 $h_{W,b}(x)$ 含义相同.
ρ	稀疏值, 可以用它指定我们所需的稀疏程度
$\hat{\rho}_i$	(sparse autoencoder中) 隐藏单元 i 的平均激活值.
β	(sparse autoencoder目标函数中) 稀疏值惩罚项的权重.

Chapter 2

矢量化编程实现

2.1 矢量化编程

当使用学习算法时，一段更快的代码通常意味着项目进展更快。例如，如果你的学习算法需要花费20分钟运行完成，这意味着你每个小时能“尝试”3个新主意。但是假如你的程序需要20个小时来运行，这意味着你一天只能“尝试”一个新主意，因为你需要花费这么长时间来等待程序的反馈。对于后者，假如你可以提升代码的效率让其只需要运行10个小时，那么你的效率差不多提升一倍。

矢量化编程是提高算法速度的一种有效方法。为了提升特定数值运算操作（如矩阵相乘、矩阵相加、矩阵-向量乘法等）的速度，数值计算和并行计算的研究人员已经努力了几十年。矢量化编程的思想就是尽量使用这些被高度优化的数值运算操作来实现我们的学习算法。

例如，假设 $x \in \Re^{n+1}$ 和 $\theta \in \Re^{n+1}$ 为向量，需要计算 $z = \theta^T x$ ，那么可以按以下方式实现（使用Matlab）：

```
1 z = 0;
2 for i=1:(n+1),
3     z = z + theta(i) * x(i);
```

```
4 end;
```

或者可以更加简单的写为:

```
1 z = theta' * x;
```

第二段程序代码不仅简单，而且运行速度更快。

通常，一个编写Matlab/Octave程序的诀窍是：

代码中尽可能避免显式的for循环。

上面的第一段代码使用了一个显式的for循环。通过不使用for循环实现相同功能，可以显著提升运行速度。对Matlab/Octave代码进行矢量化的工作很大一部分集中在避免使用for循环上，因为这可以使得Matlab/Octave更多地利用代码中的并行性，同时其解释器的计算开销更小。关于编写代码的策略，开始时你会觉得矢量化代码更难

编写、阅读和调试，但你需要在编码和调试的便捷性与运行时间之间做个权衡。因此，刚开始编写程序的时候，你可能会选择不使用太多矢量化技巧来实现你的算法，并验证它是否正确（可能只在一个小问题上验证）。在确定它正确后，你可以每次只矢量化一小段代码，并在这段代码之后暂停，以验证矢量化后的代码计算结果和之前是否相同。最后，你会有望得到一份正确的、经过调试的、矢量化且有效率的代码。一旦对矢量化常见的方法和技巧熟悉后，你

将会发现对代码进行矢量化通常并不太费劲。矢量化可以使你的代码运行的更快，而且在某些情况下，还简化了你的代码。

2.2 逻辑回归的向量化实现样例

我们想用批量梯度上升法对logistic回归分析模型进行训练，其模型

如下：

$$h_{\theta}(x) = \frac{1}{1 + \exp(-\theta^T x)}, \quad (2.1)$$

让我们遵从公开课程视频与CS229教学讲义的符号规范，设 $x_0 = 1$ ，于是 $x \in R^{n+1}$ ， $\theta \in R^{n+1}$ ， θ_0 为截距。假设我们有 m 个训练样本 $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ ，

而批量梯度上升法的更新法则是： $\theta := \theta + \alpha \nabla_{\theta} l(\theta)$ ，这里的 $l(\theta)$ 是对数似然函数， $\nabla_{\theta} l(\theta)$ 是其导函数。[注：下文的符号规范与<公开课程视频>或<教学讲义CS229：机器学习>中的相同，详细内容可以参见公开课程视频或教学讲义[<http://cs229.stanford.edu/>] 于是，我们需要如下计算梯度：

$$\nabla_{\theta} \ell(\theta) = \sum_{i=1}^m \left(y^{(i)} - h_{\theta}(x^{(i)}) \right) x_j^{(i)}. \quad (2.2)$$

我们用Matlab/Octave风格变量 x 表示输入数据构成的样本矩阵， $x(:, i)$ 代表第 i 个训练样本 $x^{(i)}$ 而 $x(j, i)$ 就代表 $x_j^{(i)}$ （译者注：第 i 个训练样本向量的第 j 个元素）。同样，用Matlab/Octave风格变量 y 表示由训练样本集合的全体类别标号所构成的行向量，则该向量的第 i 个元素 $y(i)$ 就代表上式中的 $y^{(i)} \in \{0, 1\}$ 。（注意这里跟公开课程视频及CS229的符号规范不同，矩阵 x 按列而不是按行存放输入训练样本，同样， $y \in R^{1 \times m}$ 是行向量而不是列向量。）

以下是梯度运算代码的一种实现，非常恐怖，速度极慢：

```

1 %code 1
2 grad = zeros(n+1,1);
3 for i=1:m,
4     h = sigmoid(theta'*x(:,i));
5     temp = y(i) - h;
```

```
6   for j=1:n+1,
7       grad(j) = grad(j) + temp * x(j,i);
8   end;
9 end;
```

嵌套的for循环语句使这段代码的运行非常缓慢。以下是更典型的实现方式，它对算法进行部分向量化，带来更优的执行效率：

```
1 % code2
2 grad = zeros(n+1,1);
3 for i=1:m,
4     grad = grad + (y(i) - sigmoid(theta'*x(:,i))) * x(:,i);
5 end;
```

但是，或许可以向量化得更彻底些。如果去除for循环，我们就可以显著地改善代码执行效率。特别的，假定b是一个列向量，A是一个矩阵，我们用以下两种方式来计算A*b：

```
1
2 grad = zeros(n+1,1);
3 for i=1:m,
4     grad = grad + b(i) * A(:,i);
5 end;
6
7
8 grad = A*b;
```

我们看到，代码2是用了低效的for循环语句执行梯度上升（译者注：原文是下降）运算，将b(i)看成 $(y(i) - \text{sigmoid}(\theta'x(:,i)))$ ，A看成x，我们就可以使用以下高效率的代码：

```
1 % code3
2 grad = x * (y- sigmoid(theta '*x));
```

这里我们假定Matlab/Octave的 $\text{sigmoid}(z)$ 函数接受一个向量形式的输入 z ，依次对输入向量的每个元素施行 sigmoid 函数，最后返回运算结果，因此 $\text{sigmoid}(z)$ 的输出结果是一个与 z 有相同维度的向量。当训练数据集很大时，最终的实现（译者注：代码3）充分发挥了Matlab/Octave高度优化的数值线性代数库的优势来进行矩阵-向量操作，因此，比起之前代码要高效得多。想采用向量化实现并非易事，通常需要周密的思考。但当你熟练掌握向量化操作后，你会发现，这里面有固定的设计模式（对应少量的向量化技巧），可以灵活运用到很多不同的代码片段中。

2.3 神经网络矢量化

在本节，我们将引入神经网络的向量化版本。在前面关于神经网络介绍的章节中，我们已经给出了一个部分向量化的实现，它在一次输入一个训练样本时是非常有效率的。下边我们看看如何实现同时处理多个训练样本的算法。具体来讲，我们将把正向传播、反向传播这两个步骤以及稀疏特征集学习扩展为多训练样本版本。

正向传播

考虑一个三层网络(一个输入层、一个隐含层、以及一个输出层)，并且假定 x 是包含一个单一训练样本 $x^{(i)} \in \Re^n$ 的列向量。则向量化的正向传播步骤如下：

$$z^{(2)} = W^{(1)}x + b^{(1)} \quad (2.3)$$

$$a^{(2)} = f(z^{(2)}) \quad (2.4)$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)} \quad (2.5)$$

$$h_{W,b}(x) = a^{(3)} = f(z^{(3)}) \quad (2.6)$$

这对于单一训练样本而言是非常有效的一种实现，但是当我们需要处理 m 个训练样本时，则需要把如上步骤放入一个for循环中。

更具体点来说，参照逻辑回归向量化的例子，我们用Matlab/Octave风格变量 x 表示包含输入训练样本的矩阵， $x(:, i)$ 代表第 i 个训练样本。

则 x 正向传播步骤可如下实现：

```

1 %向量化实现
2 for i=1:m,
3     z2 = W1 * x(:, i) + b1;
4     a2 = f(z2);
5     z3 = W2 * a2 + b2;
6     h(:, i) = f(z3);
7 end;
```

这个for循环能否去掉呢？对于很多算法而言，我们使用向量来表示计算过程中的中间结果。例如在前面的非向量化实现中， $z2, a2, z3$ 都是列向量，分别用来计算隐层和输出层的激励结果。为了充分利用并行化和高效矩阵运算的优势，我们希望算法能同时处理多个训练样本。让我们先暂时忽略前面公式中的 $b1$ 和 $b2$ (把它们设置为0)，那么可以实现如下：

```

1 % 向量化实现 (忽略 b1, b2)
2 z2 = W1 * x;
```



```

3 a2 = f(z2);
4 z3 = W2 * a2;
5 h = f(z3)

```

在这个实现中， z_2, a_2, z_3 都是矩阵，每个训练样本对应矩阵的一列。在对多个训练样本实现向量化时常用的设计模式是，虽然前面每个样本对应一个列向量（比如 z_2 ），但我们可把这些列向量堆叠成一个矩阵以充分享受矩阵运算带来的好处。这样，在这个例子中， a_2 就成了一个 $s_2 \times m$ 的矩阵（ s_2 是网络第二层中的神经元数， m 是训练样本个数）。矩阵 a_2 的物理含义是，当第 i 个训练样本 $x(i)$ 输入到网络中时，它的第 i 列就表示这个输入信号对隐神经元（网络第二层）的激励结果。

在上面的实现中，我们假定激活函数 $f(z)$ 接受矩阵形式的输入 z ，并对输入矩阵按列分别施以激活函数。需要注意的是，你在实现 $f(z)$ 的时候要尽量多用Matlab/Octave的矩阵操作，并尽量避免使用for循环。假定激活函数采用Sigmoid函数，则实现代码如下所示：

```

1 % 低效的、非向量化的激活函数实现
2 function output = unvectorized_f(z)
3 output = zeros(size(z))
4 for i=1:size(z,1),
5     for j=1:size(z,2),
6         output(i,j) = 1/(1+exp(-z(i,j)));
7     end;
8 end;
9 end
10 % 高效的、向量化激活函数实现
11 function output = vectorized_f(z)

```

```

12 output = 1./(1+exp(-z));      % "./" 在Matlab或Octave中表示
    对矩阵的每个元素分别进行除法操作
13 end

```

最后，我们上面的正向传播向量化实现中忽略了b1和b2，现在要把他们包含进来，为此我们需要用到Matlab/Octave的内建函数repmat：

```

1 % 正向传播的向量化实现
2 z2 = W1 * x + repmat(b1,1,m);
3 a2 = f(z2);
4 z3 = W2 * a2 + repmat(b2,1,m);
5 h = f(z3)

```

repmat(b1,1,m)的运算效果是，它把列向量b1拷贝m份，然后堆叠成如下矩阵：

$$\begin{bmatrix} | & | & & | \\ b1 & b1 & \cdots & b1_s \\ | & | & & | \end{bmatrix}.$$

这就构成一个 $s_2 \times m$ 的矩阵。它和 $W1 * x$ 相加，就等于是把 $W1 * x$ 矩阵（译者注：这里x是训练矩阵而非向量，所以 $W1 * x$ 代表两个矩阵相乘，结果还是一个矩阵）的每一列加上b1。如果不熟悉的话，可以参考Matlab/Octave的帮助文档获取更多信息（输入“help repmat”）。repmat作为Matlab/Octave的内建函数，运行起来是相当高效的，远远快过我们自己用for循环实现的效果。

反向传播

现在我们来描述反向传播向量化的思路。在阅读这一节之前，强烈建议各位仔细阅读前面介绍的正向传播的例子代码，确保你已经完全理解。下边我们只会给出反向传播向量化实现的大致纲要，而

由你来完成具体细节的推导（见向量化练习）。

对于监督学习，我们有一个包含 m 个带类别标号样本的训练集 $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ 。（对于自编码网络，我们只需令 $y(i) = x(i)$ 即可，但这里考虑的是更一般的情况。）

假定网络的输出有 s_3 维，因而每个样本的类别标号向量就记为 $y^{(i)} \in \mathbb{R}^{s_3}$ 。在我们的Matlab/Octave数据结构实现中，把这些输出按列合在一起形成一个Matlab/Octave风格变量 y ，其中第 i 列 $y(:,i)$ 就是 $y(i)$ 。

现在我们要计算梯度项 $\nabla_{W^{(l)}} J(W, b)$ 和 $\nabla_{b^{(l)}} J(W, b)$ 。对于梯度中的第一项，就像过去在反向传播算法中所描述的那样，对于每个训练样本 (x, y) ，我们可以这样来计算：

$$\delta^{(3)} = -(y - a^{(3)}) \bullet f'(z^{(3)}), \quad (2.7)$$

$$\delta^{(2)} = ((W^{(2)})^T \delta^{(3)}) \bullet f'(z^{(2)}), \quad (2.8)$$

$$\nabla_{W^{(2)}} J(W, b; x, y) = \delta^{(3)} (a^{(2)})^T, \quad (2.9)$$

$$\nabla_{W^{(1)}} J(W, b; x, y) = \delta^{(2)} (a^{(1)})^T. \quad (2.10)$$

在这里 \bullet 表示对两个向量按对应元素相乘的运算（译者注：其结果还是一个向量）。为了描述简单起见，我们这里暂时忽略对参数 $b^{(l)}$ 的求导，不过在你真正实现反向传播时，还是需要计算关于它们的导数的。

假定我们已经实现了向量化的正向传播方法，如前面那样计算了矩阵形式的变量 z_2 , a_2 , z_3 和 h ，那么反向传播的非向量化版本可如下实现：

```
1 gradW1 = zeros(size(W1));
2 gradW2 = zeros(size(W2));
3 for i=1:m,
```

```

4   delta3 = -(y(:,i) - h(:,i)) .* fprime(z3(:,i));
5   delta2 = W2'*delta3(:,i) .* fprime(z2(:,i));
6
7   gradW2 = gradW2 + delta3*a2(:,i)';
8   gradW1 = gradW1 + delta2*a1(:,i)';
9 end;

```

在这个实现中，有一个for循环。而我们想要一个能同时处理所有样本、且去除这个for循环的向量化版本。

为做到这一点，我们先把向量delta3和delta2替换为矩阵，其中每列对应一个训练样本。我们还要实现一个函数 $f_{prime}(z)$ ，该函数接受矩阵形式的输入 z ，并且对矩阵的按元素分别执行 $f'(\cdot)$ 。这样，上面for循环中的4行Matlab代码中每行都可单独向量化，以一行新的（向量化的）Matlab代码替换它（不再需要外层的for循环）。

在向量化练习中，我们要求你自己去推导出这个算法的向量化版本。如果你已经能从上面的描述中了解如何去做，那么我们强烈建议你实践一下。虽然我们已经为你准备了反向传播的向量化实现提示，但还是鼓励你在不看提示的情况下自己去推导一下。

稀疏自编码网络

稀疏自编码网络中包含一个额外的稀疏惩罚项，目的是限制神经元的平均激活率，使其接近某个（预设的）目标激活率 ρ 。其实在对单个训练样本上执行反向传播时，我们已经考虑了如何计算这个稀疏惩罚项，如下所示：

$$\delta_i^{(2)} = \left(\left(\sum_{j=1}^{s_2} W_{ji}^{(2)} \delta_j^{(3)} \right) + \beta \left(-\frac{\rho}{\hat{\rho}_i} + \frac{1-\rho}{1-\hat{\rho}_i} \right) \right) f'(z_i^{(2)}). \quad (2.11)$$

在非向量化的实现中，计算代码如下：

```
1 % 稀疏惩罚 Delta
2 sparsity_delta =
3     - rho ./ rho_hat + (1 - rho) ./ (1 - rho_hat);
4 for i=1:m,
5     ...
6     delta2 = (W2'*delta3(:,i)
7         + beta*sparsity_delta) .* fprime(z2(:,i));
8     ...
9 end;
```

但在上面的代码中，仍旧含有一个需要在整个训练集上运行的for循环，这里delta2是一个列向量。

Chapter 3

预处理：主成分分析与白化

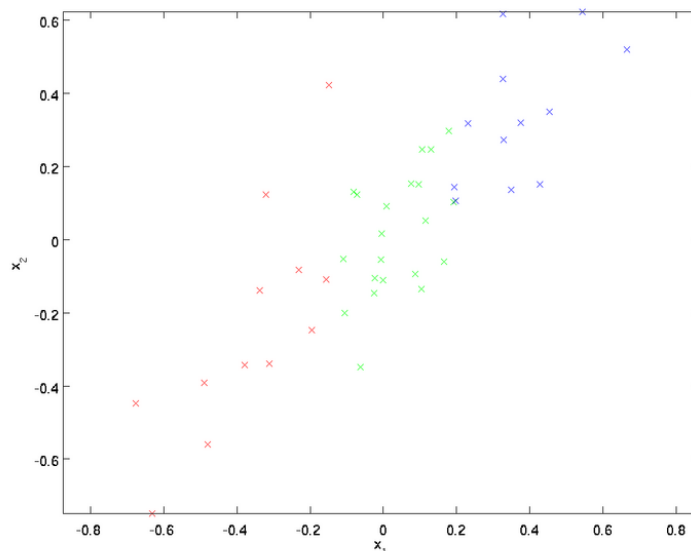
3.1 主成分分析

引言

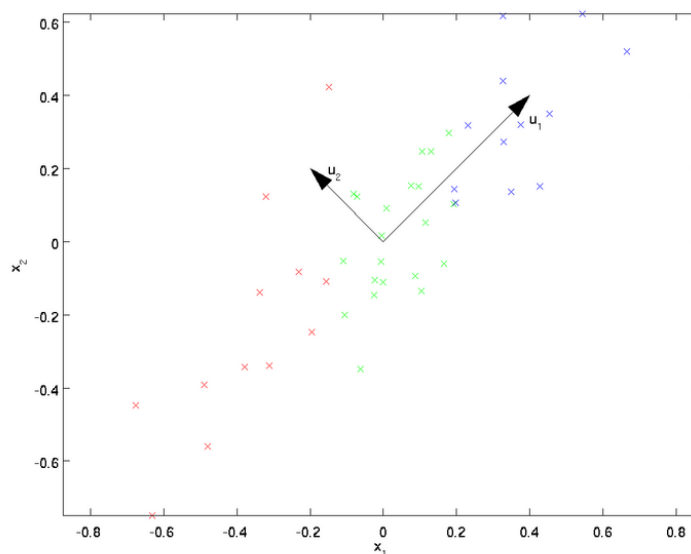
主成分分析（PCA）是一种能够极大提升无监督特征学习速度的数据降维算法。更重要的是，理解PCA算法，对实现白化算法有很大的帮助，很多算法都先用白化算法作预处理步骤。假设你使用图像来训练算法，因为图像中相邻的像素高度相关，输入数据是有一定冗余的。具体来说，假如我们正在训练的16x16灰度值图像，记为一个256维向量 $x \in \mathbb{R}^{256}$ ，其中特征值 x_j 对应每个像素的亮度值。由于相邻像素间的相关性，PCA算法可以将输入向量转换为一个维数低很多的近似向量，而且误差非常小。

实例和数学背景

在我们的实例中，使用的输入数据集表示为 $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ ，维度 $n = 2$ 即 $x^{(i)} \in \mathbb{R}^2$ 。假设我们想把数据从2维降到1维。（在实际应用中，我们也许需要把数据从256维降到50维；在这里使用低维数据，主要是为了更好地可视化算法的行为）。下图是我们的数据集：这些数据已经进行了预处理，使得每个特征 x_1 和 x_2 具有相同的均值（零）和方差。为方便展示，根据 x_1 值的大小，我们将每个点分别涂上了三种颜色之一，但该颜色并不用于算法而仅用于图解。PCA算法将寻找一个低维空间来投影我们的数据。从下图中可



以看出， u_1 是数据变化的主方向，而 u_2 是次方向。也就是说，数



据在 u_1 方向上的变化要比在 u_2 方向上大。为更形式化地找出方向 u_1 和 u_2 ，我们首先计算出矩阵 Σ ，如下所示：

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T. \quad (3.1)$$

假设 x 的均值为零，那么 Σ 就是 x 的协方差矩阵。（符号 Σ ，读 "Sigma"，是协方差矩阵的标准符号。虽然看起来与求和符号 $\sum_{i=1}^n i$ 比较像，但它们其实是两个不同的概念。）可以证明，数

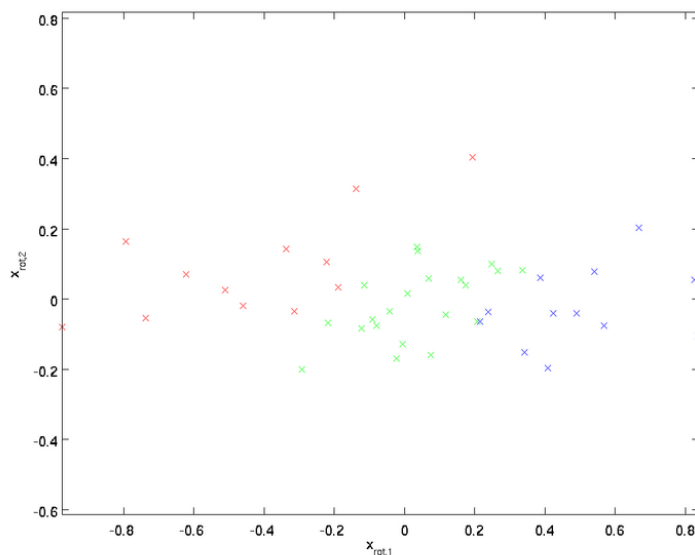
据变化的主方向 u_1 就是协方差矩阵 Σ 的主特征向量，而 u_2 是次特征向量。注：如果你对如何得到这个结果的具体数学推导过程感兴趣，可以参看CS229（机器学习）PCA部分的课件（链接在本页底部）。但如果仅仅是想跟上本课，可以不必如此。你可以通过标准的数值线性代数运算软件求得特征向量（见实现说明）。我们先计算出协方差矩阵 Σ 的特征向量，按列排放，而组成矩阵 U ：

$$U = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \cdots & u_n \\ | & | & & | \end{bmatrix} \quad (3.2)$$

旋转数据

至此，我们可以把 x 用 (u_1, u_2) 基表达为：

$$x_{\text{rot}} = U^T x = \begin{bmatrix} u_1^T x \\ u_2^T x \end{bmatrix} \quad (3.3)$$



这就是把训练数据集旋转到 u_1, u_2 基后的结果。一般而言，运算 $U^T x$ 表示旋转到基 u_1, u_2, \dots, u_n 之上的训练数据。矩阵 U 有正交性，即满足 $U^T U = U U^T = I$ ，所以若想将旋转后的向量 x_{rot} 还原为原始数据 x ，将其左乘矩阵 U 即可： $x = U x_{\text{rot}}$ ，验算一下：

$$Ux_{\text{rot}} = UU^T x = x.$$

数据降维

数据的主方向就是旋转数据的第一维 $x_{\text{rot},1}$ 。因此，若想把这数据降到一维，可令：

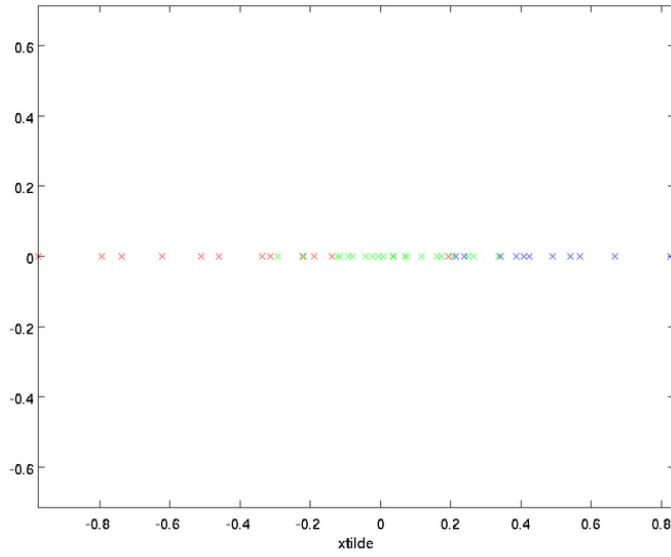
$$\tilde{x}^{(i)} = x_{\text{rot},1}^{(i)} = u_1^T x^{(i)} \in \mathbb{R}. \quad (3.4)$$

更一般的，假如想把数据 $x \in \mathbb{R}^n$ 降到 k 维表示 $\tilde{x} \in \mathbb{R}^k$ （令 $k < n$ ），只需选取 x_{rot} 的前 k 个成分，分别对应前 k 个数据变化的主方向。PCA的另外一种解释是： x_{rot} 是一个 n 维向量，其中前几个成分可能比较大（例如，上例中大部分样本第一个成分 $x_{\text{rot},1}^{(i)} = u_1^T x^{(i)}$ 的取值相对较大），而后面成分可能会比较小（例如，上例中大部分样本的 $x_{\text{rot},2}^{(i)} = u_2^T x^{(i)}$ 较小）。PCA算法做的其实就是丢弃 x_{rot} 中后面（取值较小）的成分，就是将这些成分的值近似为零。具体的说，设 \tilde{x} 是 x_{rot} 的近似表示，那么将 x_{rot} 除了前 k 个成分外，其余全赋值为零，就得到：

$$\tilde{x} = \begin{bmatrix} x_{\text{rot},1} \\ \vdots \\ x_{\text{rot},k} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \approx \begin{bmatrix} x_{\text{rot},1} \\ \vdots \\ x_{\text{rot},k} \\ x_{\text{rot},k+1} \\ \vdots \\ x_{\text{rot},n} \end{bmatrix} = x_{\text{rot}} \quad (3.5)$$

在本例中，可得 \tilde{x} 的点图如下（取 $n=2, k=1$ ）：

然而，由于上面 \tilde{x} 的后 $n-k$ 项均为零，没必要把这些零项保留下来。所以，我们仅用前 k 个（非零）成分来定义 k 维向量 \tilde{x} 。这也解释了我们为什么会以 u_1, u_2, \dots, u_n 为基来表示数据：要决定保留哪些成分变得很简单，只需取前 k 个成分即可。这时也可以说，我们“保留了前 k 个PCA（主）成分”。

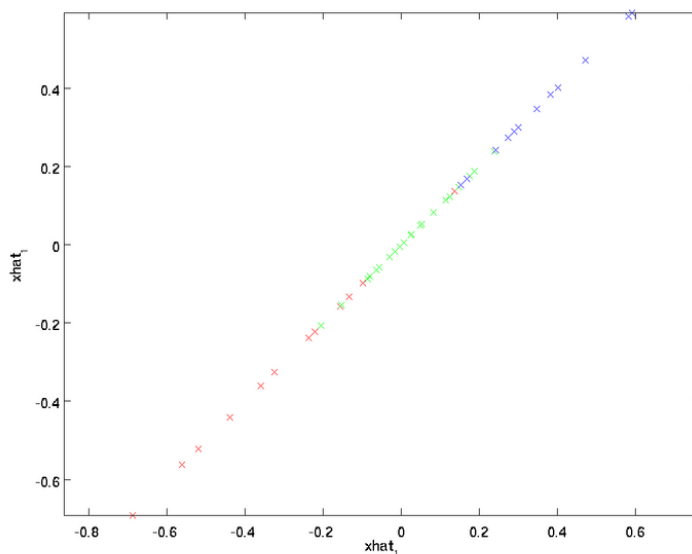


还原近似数据 现在，我们得到了原始数据 $x \in \Re^n$ 的低维“压缩”表征量 $\tilde{x} \in \Re^k$ ，反过来，如果给定 \tilde{x} ，我们应如何还原原始数据 x 呢？查看以往章节可知，要转换回来，只需 $x = Ux_{\text{rot}}$ 即可。进一步，我们把 \tilde{x} 看作将 x_{rot} 的最后 $n - k$ 个元素被置0所得的近似表示，因此如果给定 $\tilde{x} \in \Re^k$ ，可以通过在其末尾添加 $n - k$ 个0来得到对 $x_{\text{rot}} \in \Re^n$ 的近似，最后，左乘 U 便可近似还原出原数据 x 。具体来说，计算如下：

$$\hat{x} = U \begin{bmatrix} \tilde{x}_1 \\ \vdots \\ \tilde{x}_k \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \sum_{i=1}^k u_i \tilde{x}_i. \quad (3.6)$$

上面的等式基于先前对 U 的定义。在实现时，我们实际上并不先给 \tilde{x} 填0然后再左乘 U ，因为这意味着大量的乘0运算。我们可用 $\tilde{x} \in \Re^k$ 来与 U 的前 k 列相乘，即上式中最右项，来达到同样的目的。将该算法应用于本例中的数据集，可得如下关于重构数据 \hat{x} 的

点图：



由图可见，我们得到的是对原始数据集的一维近似重构。在训练自动编码器或其它无监督特征学习算法时，算法运行时间将依赖于输入数据的维数。若用 $\tilde{x} \in \mathbb{R}^k$ 取代 x 作为输入数据，那么算法就可使用低维数据进行训练，运行速度将显著加快。对于很多数据集来说，低维表征量 \tilde{x} 是原数据集的极佳近似，因此在这些场合使用PCA是很合适的，它引入的近似误差的很小，却可显著地提高你算法的运行速度。

选择主成分个数

我们该如何选择 k ，即保留多少个PCA主成分？在上面这个简单的二维实验中，保留第一个成分看起来是自然的选择。对于高维数据来说，做这个决定就没那么简单：如果 k 过大，数据压缩率不高，在极限情况 $k = n$ 时，等于是在使用原始数据（只是旋转投射到了不同的基）；相反地，如果 k 过小，那数据的近似误差太大。决定 k 值时，我们通常会考虑不同 k 值可保留的方差百分比。具体来说，如果 $k = n$ ，那么我们得到的是对数据的完美近似，也就是保留了100%的方差，即原始数据的所有变化都被保留下来；相

反，如果 $k = 0$ ，那等于是使用零向量来逼近输入数据，也就是只有0%的方差被保留下来。一般而言，设 $\lambda_1, \lambda_2, \dots, \lambda_n$ 表示 Σ 的特征值（按由大到小顺序排列），使得 λ_j 为对应于特征向量 u_j 的特征值。那么如果我们保留前 k 个成分，则保留的方差百分比可计算为：

$$\frac{\sum_{j=1}^k \lambda_j}{\sum_{j=1}^n \lambda_j}. \quad (3.7)$$

在上面简单的二维实验中， $\lambda_1 = 7.29$ ， $\lambda_2 = 0.69$ 。因此，如果保留 $k = 1$ 个主成分，等于我们保留了 $7.29/(7.29 + 0.69) = 0.913$ ，即91.3%的方差。对保留方差的百分比进行更正式的定义已超出了本教程的范围，但很容易证明， $\lambda_j = \sum_{i=1}^m x_{\text{rot},j}^2$ 。因此，如果 $\lambda_j \approx 0$ ，则说明 $x_{\text{rot},j}$ 也就基本上接近于0，所以用0来近似它并不会产生多大损失。这也解释了为什么要保留前面的主成分（对应的 λ_j 值较大）而不是末尾的那些。这些前面的主成分 $x_{\text{rot},j}$ 变化性更大，取值也更大，如果将其设为0势必引入较大的近似误差。以处理图像数据为例，一个惯常的经验法则是选择 k 以保留99%的方差，换句话说，我们选取满足以下条件的最小 k 值：

$$\frac{\sum_{j=1}^k \lambda_j}{\sum_{j=1}^n \lambda_j} \geq 0.99. \quad (3.8)$$

对其它应用，如不介意引入稍大的误差，有时也保留90-98%的方差范围。若向他人介绍PCA算法详情，告诉他们你选择的 k 保留了95%的方差，比告诉他们你保留了前120个（或任意某个数字）主成分更好理解。

对图像数据应用PCA算法

为使PCA算法能有效工作，通常我们希望所有的特征 x_1, x_2, \dots, x_n 都有相似的取值范围（并且均值接近于0）。如果你曾在其它应用中使用过PCA算法，你可能知道有必要单独对每个特征做预处理，即

通过估算每个特征 x_j 的均值和方差，而后将其取值范围规整化为零均值和单位方差。但是，对于大部分图像类型，我们却不需要进行这样的预处理。假定我们将在自然图像上训练算法，此时特征 x_j 代表的是像素 j 的值。所谓“自然图像”，不严格的说，是指人或动物在他们一生中所见的那种图像。注：通常我们选取含草木等内容的户外场景图片，然后从中随机截取小图像块（如16x16像素）来训练算法。在实践中我们发现，大多数特征学习算法对训练图片的确切类型并不敏感，所以大多数用普通照相机拍摄的图片，只要不是特别的模糊或带有非常奇怪的人工痕迹，都可以使用。在自然图像上进行训练时，对每一个像素单独估计均值和方差意义不大，因为（理论上）图像任一部分的统计性质都应该和其它部分相同，图像的这种特性被称作平稳性（stationarity）。具体而言，为使PCA算法正常工作，我们通常需要满足以下要求：(1)特征的均值大致为0；(2)不同特征的方差值彼此相似。对于自然图片，即使不进行方差归一化操作，条件(2)也自然满足，故而我们不再进行任何方差归一化操作（对音频数据,如声谱,或文本数据,如词袋向量，我们通常也不进行方差归一化）。实际上，PCA算法对输入数据具有缩放不变性，无论输入数据的值被如何放大（或缩小），返回的特征向量都不改变。更正式的说：如果将每个特征向量 x 都乘以某个正数（即所有特征量被放大或缩小相同的倍数），PCA的输出特征向量都将不会发生变化。既然我们不做方差归一化，唯一还需进行的规整化操作就是均值规整化，其目的是保证所有特征的均值都在0附近。根据应用，在大多数情况下，我们并不关注所输入图像的整体明亮程度。比如在对象识别任务中，图像的整体明亮程度并不会影响图像中存在的是什么物体。更为正式地说，我们对图像块的平均亮度值不感兴趣，所以可以减去这个值来进行均值规整化。具体的步骤是，如果 $x^{(i)} \in \mathbb{R}^n$ 代表16x16的图像块的亮度（灰

度) 值 ($n = 256$)，可用如下算法来对每幅图像进行零均值化操作： $\mu^{(i)} := \frac{1}{n} \sum_{j=1}^n x_j^{(i)} x_j^{(i)} := x_j^{(i)} - \mu^{(i)}$, for all j

请注意：1) 对每个输入图像块 $x^{(i)}$ 都要单独执行上面两个步骤，2) 这里的 $\mu^{(i)}$ 是指图像块 $x^{(i)}$ 的平均亮度值。尤其需要注意的是，这和为每个像素 x_j 单独估算均值是两个完全不同的概念。如果你处理的图像并非自然图像（比如，手写文字，或者白背景正中摆放单独物体），其他规整化操作就值得考虑了，而哪种做法最合适也取决于具体应用场合。但对自然图像而言，对每幅图像进行上述的零均值规整化，是默认而合理的处理。

3.2 白化

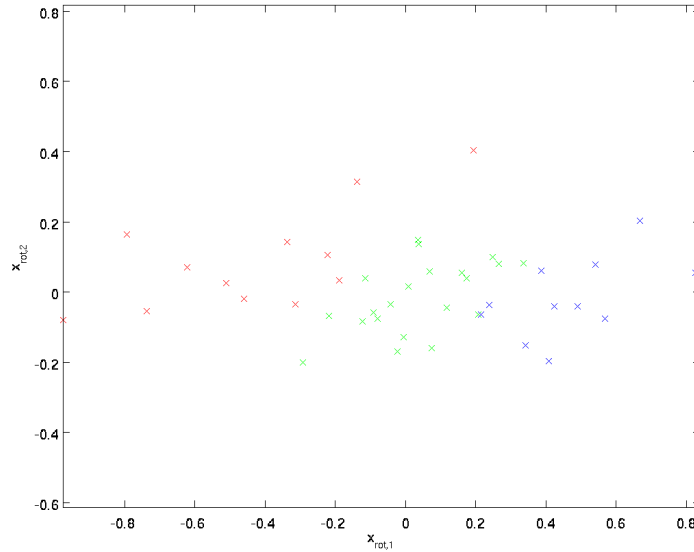
介绍

我们已经了解了如何使用PCA降低数据维度。在一些算法中还需要一个与之相关的预处理步骤，这个预处理过程称为白化（一些文献中也叫sphering）。举例来说，假设训练数据是图像，由于图像中相邻像素之间具有很强的相关性，所以用于训练时输入是冗余的。白化的目的就是降低输入的冗余性；更正式的说，我们希望通过白化过程使得学习算法的输入具有如下性质：(i)特征之间相关性较低；(ii)所有特征具有相同的方差。

2D 的例子

下面我们先从前文的2D例子描述白化的主要思想，然后分别介绍如何将白化与平滑和PCA相结合。如何消除输入特征之间的相关性？在前文计算 $x_{\text{rot}}^{(i)} = U^T x^{(i)}$ 时实际上已经消除了输入特征 $x^{(i)}$ 之间的相关性。得到的新特征 x_{rot} 的分布如下图所示：这个数据的协方差矩阵如下：

$$\begin{bmatrix} 7.29 & 0 \\ 0 & 0.69 \end{bmatrix}. \quad (3.9)$$



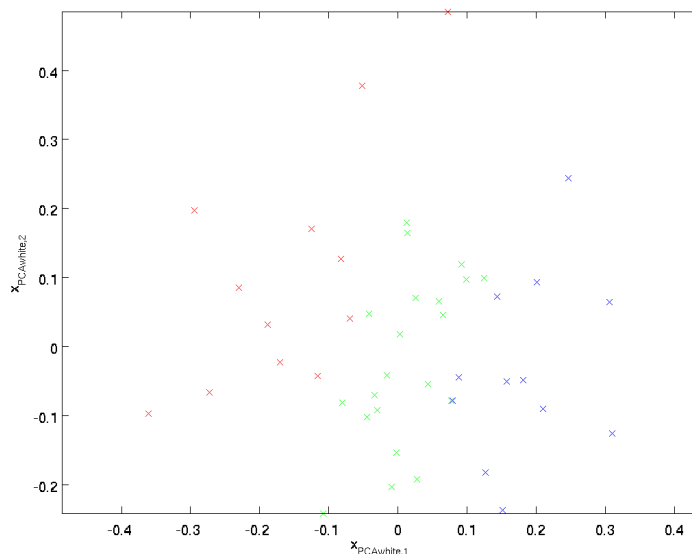
(注：严格地讲，这部分许多关于“协方差”的陈述仅当数据均值为0时成立。下文的论述都隐式地假定这一条件成立。不过即使数据均值不为0，下文的说法仍然成立，所以你无需担心这个。)

x_{rot} 协方差矩阵对角元素的值为 λ_1 和 λ_2 绝非偶然。并且非对角元素值为0；因此， $x_{\text{rot},1}$ 和 $x_{\text{rot},2}$ 是不相关的，满足我们对白化结果的第一个要求（特征间相关性降低）。为了使每个输入特征具有单位方差，我们可以直接使用 $1/\sqrt{\lambda_i}$ 作为缩放因子来缩放每个特征 $x_{\text{rot},i}$ 。具体地，我们定义白化后的数据 $x_{\text{PCAwhite}} \in \mathbb{R}^n$ 如下：

$$x_{\text{PCAwhite},i} = \frac{x_{\text{rot},i}}{\sqrt{\lambda_i}}. \quad (3.10)$$

绘制出 x_{PCAwhite} ，我们得到：

这些数据现在的协方差矩阵为单位矩阵 I 。我们说， x_{PCAwhite} 是数据经过PCA白化后的版本： x_{PCAwhite} 中不同的特征之间不相关并且具有单位方差。白化与降维相结合。如果你想要得到经过白化后的数据，并且比初始输入维数更低，可以仅保留 x_{PCAwhite} 中前 k 个成分。当我们把PCA白化和正则化结合起来时（在稍后讨论）， x_{PCAwhite} 中最后的少量成分将总是接近于0，因而舍弃这些成分不会带来很大的问题。



ZCA白化

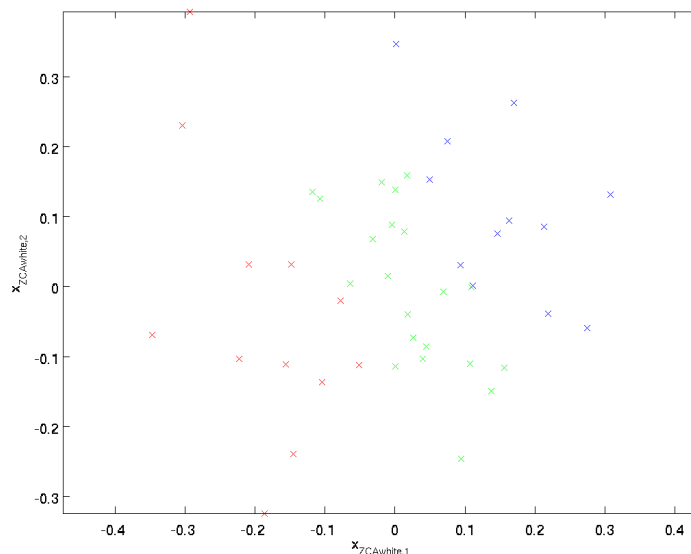
最后要说明的是，使数据的协方差矩阵变为单位矩阵 I 的方式并不唯一。具体地，如果 R 是任意正交矩阵，即满足 $RR^T = R^T R = I$ (说它正交不太严格， R 可以是旋转或反射矩阵)，那么 $R x_{\text{PCAwhite}}$ 仍然具有单位协方差。在ZCA白化中，令 $R = U$ 。我们定义ZCA白化的结果为：

$$x_{\text{ZCAwhite}} = U x_{\text{PCAwhite}} \quad (3.11)$$

绘制 x_{ZCAwhite} ，得到：

可以证明，对所有可能的 R ，这种旋转使得 x_{ZCAwhite} 尽可能地接近原始输入数据 x 。当使用 ZCA白化时(不同于 PCA白化)，我们通常保留数据的全部 n 个维度，不尝试去降低它的维数。正则化

实践中需要实现PCA白化或ZCA白化时，有时一些特征值 λ_i 在数值上接近于0，这样在缩放步骤时我们除以 $\sqrt{\lambda_i}$ 将导致除以一个接近0的值；这可能使数据上溢(赋为大数值)或造成数值不稳定。因而在实践中，我们使用少量的正则化实现这个缩放过程，即在取平



方根和倒数之前给特征值加上一个很小的常数 ϵ :

$$x_{\text{PCAwhite},i} = \frac{x_{\text{rot},i}}{\sqrt{\lambda_i + \epsilon}}. \quad (3.12)$$

当 x 在区间 $[-1, 1]$ 上时, 一般取值为 $\epsilon \approx 10^{-5}$ 。对图像来说, 这里加上 ϵ , 对输入图像也有一些平滑(或低通滤波)的作用。这样处理还能消除在图像的像素信息获取过程中产生的噪声, 改善学习到的特征(细节超出了本文的范围)。

ZCA 白化是一种数据预处理方法, 它将数据从 x 映射到 x_{ZCAwhite} 。事实证明这也是一种生物眼睛(视网膜)处理图像的粗糙模型。具体而言, 当你的眼睛感知图像时, 由于一幅图像中相邻的部分在亮度上十分相关, 大多数临近的“像素”在眼中被感知为相近的值。因此, 如果人眼需要分别传输每个像素值(通过视觉神经)到大脑中, 会非常不划算。取而代之的是, 视网膜进行一个与ZCA中相似的去相关操作(这是由视网膜上的ON-型和OFF-型光感受器细胞将光信号转变为神经信号完成的)。由此得到对输入图像的更低冗余的表示, 并将它传输到大脑。

3.3 实现主成分分析和白化

在这一节里，我们将总结PCA, PCA白化和ZCA白化算法，并描述如何使用高效的线性代数库来实现它们。首先，我们需要确保数据的均值（近似）为零。对于自然图像，我们通过减去每个图像块(patch)的均值（近似地）来达到这一目标。为此，我们计算每个图像块的均值，并从每个图像块中减去它的均值。（译注：参见PCA一章中“对图像数据应用PCA算法”一节）。Matlab实现如下：

```
1 avg = mean(x, 1); % 分别为每个图像块计算像素强度的均值。
2 x = x - repmat(avg, size(x, 1), 1);
```

下面，我们要计算 $\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T$ ，如果你在Matlab中实现（或者在C++, Java等中实现，但可以使用高效的线性代数库），直接求和效率很低。不过，我们可以这样一气呵成。

```
1 sigma = x * x' / size(x, 2);
```

（自己推导一下看看）这里，我们假设 x 为一数据结构，其中每列表示一个训练样本（所以 x 是一个 nm 的矩阵）。接下来，PCA计算 Σ 的特征向量。你可以使用Matlab的 `eig` 函数来计算。但是由于 Σ 是对称半正定的矩阵，用 `svd` 函数在数值计算上更加稳定。具体来说，如果你使用

```
1 [U,S,V] = svd(sigma);
```

那矩阵 U 将包含 Σ 的特征向量（一个特征向量一行，从主向量开始排序），矩阵 S 对角线上的元素将包含对应的特征值（同样降序排列）。矩阵 V 等于 U 的转置，可以忽略。（注意 `svd` 函数实际上计算的是一个矩阵的奇异值和奇异向量，就对称半正定矩阵的特殊情况来说，它们对应于特征值和特征向量，这里我们也只

关心这一特例。关于奇异向量和特征向量的详细讨论超出了本文范围。) 最后, 我们可以这样计算 x_{rot} 和 \tilde{x} :

```
1 xRot = U' * x; % 数据旋转后的结果。
2 xTilde = U(:, 1:k)' * x; % 数据降维后的结果, 这里k希望保留的特征向量的数目。
```

这以 $\tilde{x} \in \mathbb{R}^k$ 的形式给出了数据的PCA表示。顺便说一下, 如果 x 是一个包括所有训练数据的 nm 矩阵, 这也是一种向量化的实现方式, 上面的式子可以让你一次对所有的训练样本计算出 x_{rot} 和 \tilde{x} 。得到的 x_{rot} 和 \tilde{x} 中, 每列对应一个训练样本。为计算PCA白化后的数据 x_{PCAwhite} , 可以用 $x_{\text{PCAwhite}} = \text{diag}(1./\text{sqrt}(\text{diag}(S) + \text{epsilon})) * U' * x$; 因为 S 的对角线包括了特征值 λ_i , 这其实就是同时为所有样本 i 计算 $x_{\text{PCAwhite}, i} = \frac{x_{\text{rot}, i}}{\sqrt{\lambda_i}}$ 的简洁表达。最后, 你也可以这样计算ZCA白化后的数据 x_{ZCAwhite} :

```
1 xZCAwhite = U * diag(1./sqrt(diag(S) + epsilon)) * U' * x;
```

Chapter 4

Softmax回归

4.1 Softmax回归

简介

在本节中，我们介绍Softmax回归模型，该模型是logistic回归模型在多分类问题上的推广，在多分类问题中，类标签 y 可以取两个以上的值。Softmax回归模型对于诸如MNIST手写数字分类等问题是很有用的，该问题的目的是辨识10个不同的单个数字。Softmax回归是有监督的，不过后面也会介绍它与深度学习/无监督学习方法的结合。（译者注：MNIST 是一个手写数字识别库，由NYU 的Yann LeCun 等人维护。<http://yann.lecun.com/exdb/mnist/>）

回想一下在 logistic 回归中，我们的训练集由 m 个已标记的样本构成： $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ ，其中输入特征 $x^{(i)} \in \mathbb{R}^{n+1}$ 。（我们对符号的约定如下：特征向量 x 的维度为 $n+1$ ，其中 $x_0 = 1$ 对应截距项。）由于 logistic 回归是针对二分类问题的，因此类标记 $y^{(i)} \in \{0, 1\}$ 。假设函数(hypothesis function) 如下：

$$h_{\theta}(x) = \frac{1}{1 + \exp(-\theta^T x)}, \quad (4.1)$$

我们将训练模型参数 θ ，使其能够最小化代价函数：

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] \quad (4.2)$$

在 softmax回归中，我们解决的是多分类问题（相对于 logistic回归解决的二分类问题），类标 y 可以取 k 个不同的值（而不是 2 个）。因此，对于训练集 $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ ，我们有 $y^{(i)} \in \{1, 2, \dots, k\}$ 。（注意此处的类别下标从 1 开始，而不是 0）。例如，在 MNIST 数字识别任务中，我们有 $k = 10$ 个不同的类别。

对于给定的测试输入 x ，我们想用假设函数针对每一个类别 j 估算出概率值 $p(y = j|x)$ 。也就是说，我们想估计 x 的每一种分类结果出现的概率。因此，我们的假设函数将要输出一个 k 维的向量（向量元素的和为 1）来表示这 k 个估计的概率值。具体地说，我们的假设函数 $h_{\theta}(x)$ 形式如下：

$$h_{\theta}(x^{(i)}) = \begin{bmatrix} p(y^{(i)} = 1|x^{(i)}; \theta) \\ p(y^{(i)} = 2|x^{(i)}; \theta) \\ \vdots \\ p(y^{(i)} = k|x^{(i)}; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \begin{bmatrix} e^{\theta_1^T x^{(i)}} \\ e^{\theta_2^T x^{(i)}} \\ \vdots \\ e^{\theta_k^T x^{(i)}} \end{bmatrix} \quad (4.3)$$

其中 $\theta_1, \theta_2, \dots, \theta_k \in \mathbb{R}^{n+1}$ 是模型的参数。请注意 $\frac{1}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}}$ 这一项对概率分布进行归一化，使得所有概率之和为 1。

为了方便起见，我们同样使用符号 θ 来表示全部的模型参数。在实现 Softmax 回归时，将 θ 用一个 $k \times (n+1)$ 的矩阵来表示会很方便，该矩阵是将 $\theta_1, \theta_2, \dots, \theta_k$ 按行罗列起来得到的，如下所示：

$$\theta = \begin{bmatrix} -\theta_1^T - \\ -\theta_2^T - \\ \vdots \\ -\theta_k^T - \end{bmatrix} \quad (4.4)$$

代价函数

现在我们来介绍 softmax 回归算法的代价函数。在下面的公式中， $1\{\cdot\}$ 是示性函数，其取值规则为： $1\{<:O h \mathfrak{m}\} = 1$ ， $1\{<:G h \mathfrak{m}\} = 0$ 。举例来说，表达式 $1\{2 + 2 = 4\}$ 的值为1， $1\{1 + 1 = 5\}$ 的值为 0。我们的代价函数为：

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{j=1}^k 1\{y^{(i)} = j\} \log \frac{e^{\theta_j^T x^{(i)}}}{\sum_{l=1}^k e^{\theta_l^T x^{(i)}}} \right] \quad (4.5)$$

值得注意的是，上述公式是logistic回归代价函数的推广。logistic回归代价函数可以改为：

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) + y^{(i)} \log h_{\theta}(x^{(i)}) \right] \quad (4.6)$$

$$= -\frac{1}{m} \left[\sum_{i=1}^m \sum_{j=0}^1 1\{y^{(i)} = j\} \log p(y^{(i)} = j | x^{(i)}; \theta) \right] \quad (4.7)$$

可以看到，Softmax代价函数与logistic 代价函数在形式上非常类似，只是在Softmax损失函数中对类标记的 k 个可能值进行了累加。注意在Softmax回归中将 x 分类为类别 j 的概率为：

$$p(y^{(i)} = j | x^{(i)}; \theta) = \frac{e^{\theta_j^T x^{(i)}}}{\sum_{l=1}^k e^{\theta_l^T x^{(i)}} .$$

对于 $J(\theta)$ 的最小化问题，目前还没有闭式解法。因此，我们使用迭代的优化算法（例如梯度下降法，或 L-BFGS）。经过求导，我们得到梯度公式如下：

$$\nabla_{\theta_j} J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[x^{(i)} \left(1\{y^{(i)} = j\} - p(y^{(i)} = j | x^{(i)}; \theta) \right) \right] \quad (4.8)$$

让我们来回顾一下符号 " ∇_{θ_j} " 的含义。 $\nabla_{\theta_j} J(\theta)$ 本身是一个向量，它的第 l 个元素 $\frac{\partial J(\theta)}{\partial \theta_{jl}}$ 是 $J(\theta)$ 对 θ_j 的第 l 个分量的偏导数。

有了上面的偏导数公式以后，我们就可以将它代入到梯度下降法等算法中，来最小化 $J(\theta)$ 。例如，在梯度下降法的标准实现中，每一次迭代需要进行如下更新： $\theta_j := \theta_j - \alpha \nabla_{\theta_j} J(\theta) (j = 1, \dots, k)$ 。当实现 softmax 回归算法时，我们通常会使用上述代价函数的一个改进版本。具体来说，就是和权重衰减(weight decay)一起使用。我们接下来介绍使用它的动机和细节。

Softmax回归模型参数化的特点

Softmax 回归有一个不寻常的特点：它有一个“冗余”的参数集。为了便于阐述这一特点，假设我们从参数向量 θ_j 中减去了向量 ψ ，这时，每一个 θ_j 都变成了 $\theta_j - \psi (j = 1, \dots, k)$ 。此时假设函数变成了以下的式子：

$$p(y^{(i)} = j | x^{(i)}; \theta) = \frac{e^{(\theta_j - \psi)^T x^{(i)}}}{\sum_{l=1}^k e^{(\theta_l - \psi)^T x^{(i)}}} \quad (4.9)$$

$$= \frac{e^{\theta_j^T x^{(i)}} e^{-\psi^T x^{(i)}}}{\sum_{l=1}^k e^{\theta_l^T x^{(i)}} e^{-\psi^T x^{(i)}}} \quad (4.10)$$

$$= \frac{e^{\theta_j^T x^{(i)}}}{\sum_{l=1}^k e^{\theta_l^T x^{(i)}}}. \quad (4.11)$$

换句话说，从 θ_j 中减去 ψ 完全不影响假设函数的预测结果！这表明前面的 softmax 回归模型中存在冗余的参数。更正式一点来说，Softmax 模型被过度参数化了。对于任意一个用于拟合数据的假设函数，可以求出多组参数值，这些参数得到的是完全相同的假设函数 h_θ 。

进一步而言，如果参数 $(\theta_1, \theta_2, \dots, \theta_k)$ 是代价函数 $J(\theta)$ 的极小值点，那么 $(\theta_1 - \psi, \theta_2 - \psi, \dots, \theta_k - \psi)$ 同样也是它的极小值点，其中 ψ 可以为任意向量。因此使 $J(\theta)$ 最小化的解不是唯一的。（有趣的是，由于 $J(\theta)$ 仍然是一个凸函数，因此梯度下降时不会遇到局部最优解的问题。但是 Hessian 矩阵是奇异的/不可逆的，这会直接导致

采用牛顿法优化就遇到数值计算的问题)

注意，当 $\psi = \theta_1$ 时，我们总是可以将 θ_1 替换为 $\theta_1 - \psi = \vec{0}$ （即替换为全零向量），并且这种变换不会影响假设函数。因此我们可以去掉参数向量 θ_1 （或者其他 θ_j 中的任意一个）而不影响假设函数的表达能力。实际上，与其优化全部的 $k \times (n + 1)$ 个参数 $(\theta_1, \theta_2, \dots, \theta_k)$ （其中 $\theta_j \in \mathbb{R}^{n+1}$ ），我们可以令 $\theta_1 = \vec{0}$ ，只优化剩余的 $(k - 1) \times (n + 1)$ 个参数，这样算法依然能够正常工作。

在实际应用中，为了使算法实现更简单清楚，往往保留所有参数 $(\theta_1, \theta_2, \dots, \theta_n)$ ，而不任意地将某一参数设置为 0。但此时我们需要对代价函数做一个改动：加入权重衰减。权重衰减可以解决 softmax 回归的参数冗余所带来的数值问题。

权重衰减

我们通过添加一个权重衰减项 $\frac{\lambda}{2} \sum_{i=1}^k \sum_{j=0}^n \theta_{ij}^2$ 来修改代价函数，这个衰减项会惩罚过大的参数值，现在我们的代价函数变为：

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{j=1}^k 1 \{y^{(i)} = j\} \log \frac{e^{\theta_j^T x^{(i)}}}{\sum_{l=1}^k e^{\theta_l^T x^{(i)}}} \right] + \frac{\lambda}{2} \sum_{i=1}^k \sum_{j=0}^n \theta_{ij}^2 \quad (4.12)$$

有了这个权重衰减项以后 ($\lambda > 0$)，代价函数就变成了严格的凸函数，这样就可以保证得到唯一的解了。此时的 Hessian 矩阵变为可逆矩阵，并且因为 $J(\theta)$ 是凸函数，梯度下降法和 L-BFGS 等算法可以保证收敛到全局最优解。

为了使用优化算法，我们需要求得这个新函数 $J(\theta)$ 的导数，如下：

$$\nabla_{\theta_j} J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[x^{(i)} (1 \{y^{(i)} = j\} - p(y^{(i)} = j | x^{(i)}; \theta)) \right] + \lambda \theta_j \quad (4.13)$$

通过最小化 $J(\theta)$ ，我们就能实现一个可用的 softmax 回归模型。

Softmax回归与Logistic 回归的关系

当类别数 $k = 2$ 时，softmax 回归退化为 logistic 回归。这表明 softmax 回归是 logistic 回归的一般形式。具体地说，当 $k = 2$ 时，softmax 回归的假设函数为：

$$h_{\theta}(x) = \frac{1}{e^{\theta_1^T x} + e^{\theta_2^T x}} \begin{bmatrix} e^{\theta_1^T x} \\ e^{\theta_2^T x} \end{bmatrix} \quad (4.14)$$

利用softmax回归参数冗余的特点，我们令 $\psi = \theta_1$ ，并且从两个参数向量中都减去向量 θ_1 ，得到：

$$h(x) = \frac{1}{e^{\vec{0}^T x} + e^{(\theta_2 - \theta_1)^T x}} \begin{bmatrix} e^{\vec{0}^T x} \\ e^{(\theta_2 - \theta_1)^T x} \end{bmatrix} \quad (4.15)$$

$$= \begin{bmatrix} \frac{1}{1 + e^{(\theta_2 - \theta_1)^T x}} \\ \frac{e^{(\theta_2 - \theta_1)^T x}}{1 + e^{(\theta_2 - \theta_1)^T x}} \end{bmatrix} \quad (4.16)$$

$$= \begin{bmatrix} \frac{1}{1 + e^{(\theta_2 - \theta_1)^T x}} \\ 1 - \frac{1}{1 + e^{(\theta_2 - \theta_1)^T x}} \end{bmatrix} \quad (4.17)$$

因此，用 $\theta' = \theta_2 - \theta_1$ ，我们就会发现 softmax 回归器预测其中一个类别的概率为 $\frac{1}{1 + e^{(\theta')^T x}}$ ，另一个类别概率的为 $1 - \frac{1}{1 + e^{(\theta')^T x}}$ ，这与 logistic 回归是一致的。

Softmax 回归 vs. k 个二元分类器

如果你在开发一个音乐分类的应用，需要对k种类型的音乐进行识别，那么是选择使用 softmax 分类器呢，还是使用 logistic 回归算法建立 k 个独立的二元分类器呢？这一选择取决于你的类别之间是否互斥，例如，如果你有四个类别的音乐，分别为：古典音乐、乡村音乐、摇滚乐和爵士乐，那么你可以假设每个训练样本只会被打

上一个标签（即：一首歌只能属于这四种音乐类型的其中一种），此时你应该使用类别数 $k = 4$ 的softmax回归。（如果在你的数据集中，有的歌曲不属于以上四类的其中任何一类，那么你可以添加一个“其他类”，并将类别数 k 设为5。）如果你的四个类别如下：人声音乐、舞曲、影视原声、流行歌曲，那么这些类别之间并不是互斥的。例如：一首歌曲可以来源于影视原声，同时也包含人声。这种情况下，使用4个二分类的 logistic 回归分类器更为合适。这样，对于每个新的音乐作品，我们的算法可以分别判断它是否属于各个类别。现在我们来查看一个计算视觉领域的例子，你的任务是将图像分到三个不同类别中。(i) 假设这三个类别分别是：室内场景、户外城区场景、户外荒野场景。你会使用softmax回归还是 3个logistic回归分类器呢？(ii) 现在假设这三个类别分别是室内场景、黑白图片、包含人物的图片，你又会选择 softmax 回归还是多个 logistic 回归分类器呢？在第一个例子中，三个类别是互斥的，因此更适于选择softmax回归分类器。而在第二个例子中，建立三个独立的 logistic 回归分类器更加合适。

Chapter 5

自我学习与无监督特征学习

5.1 自我学习

综述

如果已经有一个足够强大的机器学习算法，为了获得更好的性能，最靠谱的方法之一是给这个算法以更多的数据。机器学习界甚至有个说法：“有时候胜出者并非有最好的算法，而是有更多的数据。”

人们总是可以尝试获取更多的已标注数据，但是这样做成本往往很高。例如研究人员已经花了相当的精力在使用类似 AMT(Amazon Mechanical Turk) 这样的工具上，以期获取更大的训练数据集。相比大量研究人员通过手工方式构建特征，用众包的方式让多人手工标注数据是一个进步，但是我们可以做得更好。具体的说，如果算法能够从未标注数据中学习，那么我们就可以轻易地获取大量无标注数据，并从中学习。自学习和无监督特征学习就是这种的算法。尽管一个单一的未标注样本蕴含的信息比一个已标注的样本要少，但是如果获取大量无标注数据（比如从互联网上下载随机的、无标注的图像、音频剪辑或者是文本），并且算法能够有效的利用它们，那么相比大规模的手工构建特征和标注数据，算法将会取得更好的性能。

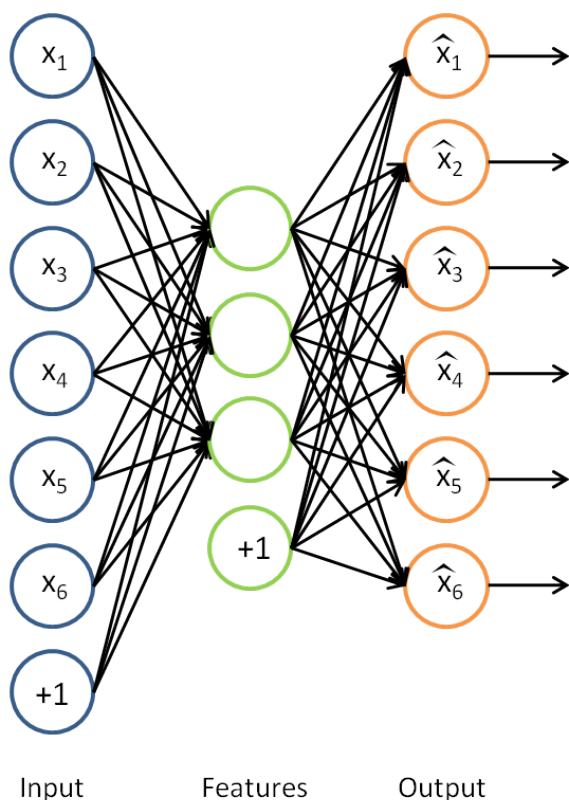
在自学习和无监督特征学习问题上，可以给算法以大量的未标注

数据，学习出较好的特征描述。在尝试解决一个具体的分类问题时，可以基于这些学习出的特征描述和任意的（可能比较少的）已标注数据，使用有监督学习方法完成分类。

在一些拥有大量未标注数据和少量的已标注数据的场景中，上述思想可能是最有效的。即使在只有已标注数据的情况下（这时我们通常忽略训练数据的类标号进行特征学习），以上想法也能得到很好的结果。

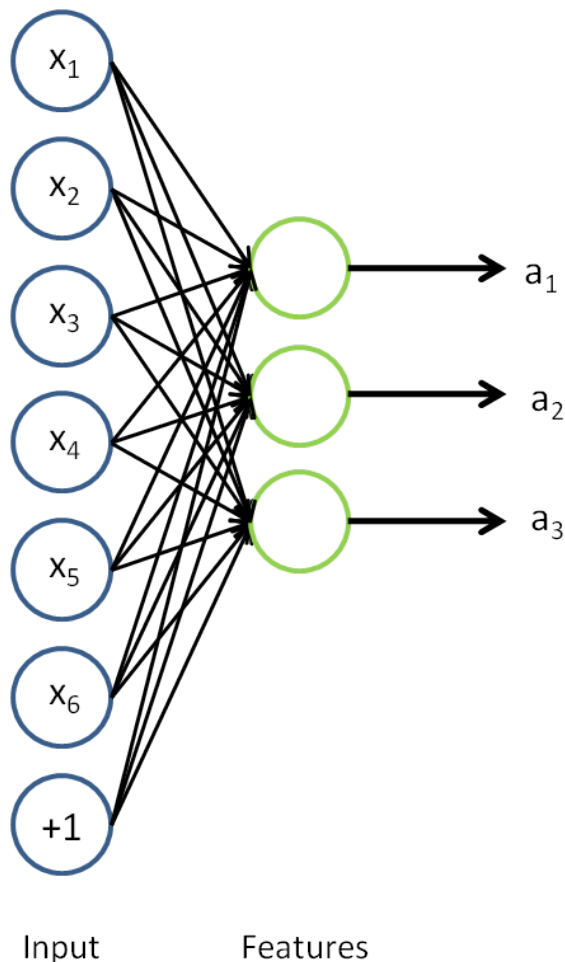
特征学习

我们已经了解到如何使用一个自编码器（autoencoder）从无标注数据中学习特征。具体来说，假定有一个无标注的训练数据集 $\{x_u^{(1)}, x_u^{(2)}, \dots, x_u^{(m_u)}\}$ （下标 u 代表“不带类标”）。现在用它们训练一个稀疏自编码器（可能需要首先对这些数据做白化或其它适当的预处理）。



利用训练得到的模型参数 $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}$ ，给定任意的输入数

据 x ，可以计算隐藏单元的激活量（activations） a 。如前所述，相比原始输入 x 来说， a 可能是一个更好的特征描述。下图的神经网络描述了特征（激活量 a ）的计算。



这实际上就是之前得到的稀疏自编码器，在这里去掉了最后一层。

假定有大小为 m_l 的已标注训练集 $\{(x_l^{(1)}, y^{(1)}), (x_l^{(2)}, y^{(2)}), \dots, (x_l^{(m_l)}, y^{(m_l)})\}$ （下标 l 表示“带类标”），我们可以为输入数据找到更好的特征描述。例如，可以将 $x_l^{(1)}$ 输入到稀疏自编码器，得到隐藏单元激活量 $a_l^{(1)}$ 。接下来，可以直接使用 $a_l^{(1)}$ 来代替原始数据 $x_l^{(1)}$ （“替代表示”，Replacement Representation）。也可以合二为一，使用新的向量 $(x_l^{(1)}, a_l^{(1)})$ 来代替原始数据 $x_l^{(1)}$ （“级联表示”，Concatenation Representation）。

经过变换后，训练集就变成 $\{(a_l^{(1)}, y^{(1)}), (a_l^{(2)}, y^{(2)}), \dots, (a_l^{(m_l)}, y^{(m_l)})\}$ 或者是 $\{((x_l^{(1)}, a_l^{(1)}), y^{(1)}), ((x_l^{(2)}, a_l^{(1)}), y^{(2)}), \dots, ((x_l^{(m_l)}, a_l^{(1)}), y^{(m_l)})\}$ （取决于使用 $a_l^{(1)}$ 替换 $x_l^{(1)}$ 还是将二者合并）。在实践中，将 $a_l^{(1)}$ 和 $x_l^{(1)}$ 合并通常表现的更好。但是考虑到内存和计算的成本，也可以使用替换操作。

最终，可以训练出一个有监督学习算法（例如 svm, logistic regression 等），得到一个判别函数对 y 值进行预测。预测过程如下：给定一个测试样本 x_{test} ，重复之前的过程，将其送入稀疏自编码器，得到 a_{test} 。然后将 a_{test} （或者 $(x_{\text{test}}, a_{\text{test}})$ ）送入分类器中，得到预测值。

数据预处理

在特征学习阶段，我们从未标注训练集 $\{x_u^{(1)}, x_u^{(2)}, \dots, x_u^{(m_u)}\}$ 中学习，这一过程中可能计算了各种数据预处理参数。例如计算数据均值并且对数据做均值标准化（mean normalization）；或者对原始数据做主成分分析（PCA），然后将原始数据表示为 $U^T x$ （又或者使用 PCA 白化或 ZCA 白化）。这样的话，有必要将这些参数保存起来，并且在后面的训练和测试阶段使用同样的参数，以保证数据进入稀疏自编码神经网络之前经过了同样的变换。例如，如果对未标注数据集进行 PCA 预处理，就必须将得到的矩阵 U 保存起来，并且应用到有标注训练集和测试集上；而不能使用有标注训练集重新估计出一个不同的矩阵 U （也不能重新计算均值并做均值标准化），否则的话可能得到一个完全不一致的数据预处理操作，导致进入自编码器的数据分布迥异于训练自编码器时的数据分布。

无监督特征学习的术语

有两种常见的无监督特征学习方式，区别在于你有什么样的未标注数据。自学习(self-taught learning) 是其中更为一般的、更强大的

学习方式，它不要求未标注数据 x_u 和已标注数据 x_l 来自同样的分布。另外一种带限制性的方式也被称为半监督学习，它要求 x_u 和 x_l 服从同样的分布。下面通过例子解释二者的区别。

假定有一个计算机视觉方面的任务，目标是区分汽车和摩托车图像；也即训练样本里面要么是汽车的图像，要么是摩托车的图像。哪里可以获取大量的未标注数据呢？最简单的方式可能是从互联网上下载一些随机的图像数据集，在这些数据上训练出一个稀疏自编码器，从中得到有用的特征。这个例子里，未标注数据完全来自于一个和已标注数据不同的分布（未标注数据集中，或许其中一些图像包含汽车或者摩托车，但是不是所有的图像都如此）。这种情形被称为自学习。

相反，如果有大量的未标注图像数据，要么是汽车图像，要么是摩托车图像，仅仅是缺失了类标号（没有标注每张图片到底是汽车还是摩托车）。也可以用这些未标注数据来学习特征。这种方式，即要求未标注样本和带标注样本服从相同的分布，有时候被称为半监督学习。在实践中，常常无法找到满足这种要求的未标注数据（到哪里找到一个每张图像不是汽车就是摩托车，只是丢失了类标号的图像数据库？）因此，自学习在无标注数据集的特征学习中应用更广。

Chapter 6

建立分类用深度网络

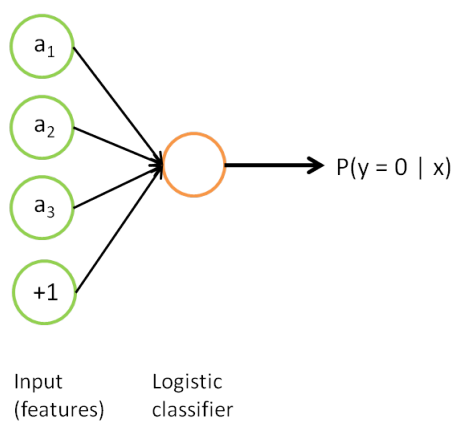
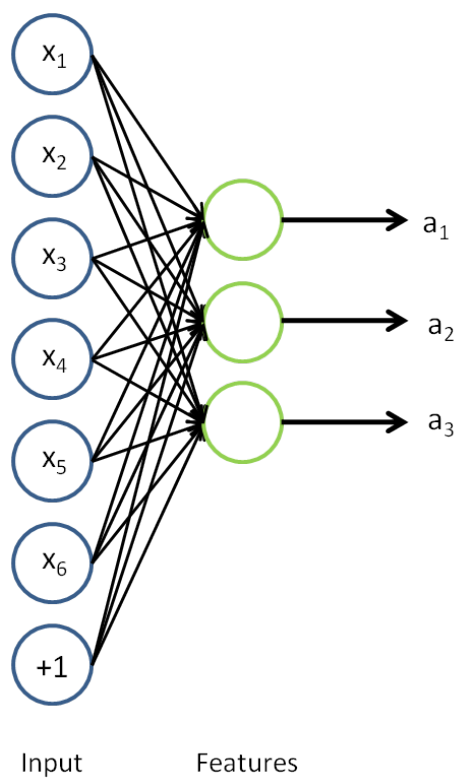
6.1 从自我学习到深层网络

在前一节中，我们利用自编码器来学习输入至 softmax 或 logistic 回归分类器的特征。这些特征仅利用未标注数据学习获得。在本节中，我们描述如何利用已标注数据进行微调，从而进一步优化这些特征。如果有大量已标注数据，通过微调就可以显著提升分类器的性能。

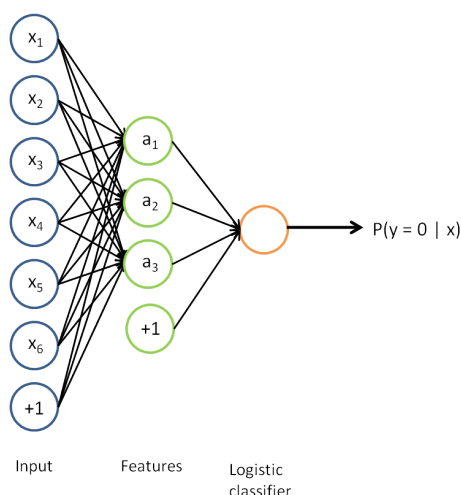
在自我学习中，我们首先利用未标注数据训练一个稀疏自编码器。随后，给定一个新样本 x ，我们通过隐含层提取出特征 a 。上述过程图示如下：

我们感兴趣的是分类问题，目标是预测样本的类别标号 y 。我们拥有标注数据集 $\{(x_l^{(1)}, y^{(1)}), (x_l^{(2)}, y^{(2)}), \dots, (x_l^{(m_l)}, y^{(m_l)})\}$ ，包含 m_l 个标注样本。此前我们已经说明，可以利用稀疏自编码器获得的特征 $a^{(l)}$ 来替代原始特征。这样就可获得训练数据集 $\{(a^{(1)}, y^{(1)}), \dots, (a^{(m_l)}, y^{(m_l)})\}$ 。最终，我们训练出一个从特征 $a^{(i)}$ 到类标号 $y^{(i)}$ 的 logistic 分类器。为说明这一过程，我们按照神经网络一节中的方式，用下图描述 logistic 回归单元（橘黄色）。

考虑利用这个方法所学到的分类器（输入-输出映射）。它描述了一个把测试样本 x 映射到预测值 $p(y = 1|x)$ 的函数。将此前的两张图片结合起来，就得到该函数的图形表示。也即，最终的分类器可以



表示为：



该模型的参数通过两个步骤训练获得：在该网络的第一层，将输入 x 映射至隐藏单元激活量 a 的权值 $W^{(1)}$ 可以通过稀疏自编码器训练过程获得。在第二层，将隐藏单元 a 映射至输出 y 的权值 $W^{(2)}$ 可以通过 logistic 回归或 softmax 回归训练获得。

这个最终分类器整体上显然是一个大的神经网络。因此，在训练获得模型最初参数（利用自动编码器训练第一层，利用 logistic/softmax 回归训练第二层）之后，我们可以进一步修正模型参数，进而降低训练误差。具体来说，我们可以对参数进行微调，在现有参数的基础上采用梯度下降或者 L-BFGS 来降低已标注样本集 $\{(x_l^{(1)}, y^{(1)}), (x_l^{(2)}, y^{(2)}), \dots, (x_l^{(m_l)}, y^{(m_l)})\}$ 上的训练误差。

使用微调时，初始的非监督特征学习步骤（也就是自动编码器和 logistic 分类器训练）有时候被称为预训练。微调的作用在于，已标注数据集也可以用来修正权值 $W^{(1)}$ ，这样可以对隐藏单元所提取的特征 a 做进一步调整。

到现在为止，我们描述上述过程时，都假设采用了“替代 (Replacement)”表示而不是“级联 (Concatenation)”表示。在替代表示中，logistic 分类器所看到的训练样本格式为 $(a^{(i)}, y^{(i)})$ ；而在级联表示中，分类器所看到的训练样本格式为 $((x^{(i)}, a^{(i)}), y^{(i)})$ 。对级

联表示同样可以进行微调（在级联表示神经网络中，输入值 x_i 也直接被输入至 logistic 分类器。对此前的神经网络示意图稍加更改，即可获得其示意图。具体的说，第一层的输入节点除了与隐层联接之外，还将越过隐层，与第三层输出节点直接相连）。但是对于微调来说，级联表示相对于替代表示几乎没有优势。因此，如果需要开展微调，我们通常使用替代表示的网络（但是如果不开展微调，级联表示的效果有时候会好得多）。

在什么时候应用微调？通常仅在有大量已标注训练数据的情况下使用。在这样的情况下，微调能显著提升分类器性能。然而，如果有大量未标注数据集（用于非监督特征学习/预训练），却只有相对较少的已标注训练集，微调的作用非常有限。

中英文对照

自我学习 self-taught learning

深层网络 deep networks

微调 fine-tune

稀疏自编码器 sparse autoencoder

梯度下降 gradient descent

非监督特征学习 unsupervised feature learning

预训练 pre-training

6.2 深度网络概览

概述

在之前的章节中，你已经构建了一个包括输入层、隐藏层以及输出层的三层神经网络。虽然该网络对于MNIST手写数字数据库非常有效，但是它还是一个非常“浅”的网络。这里的“浅”指的是特征（隐藏层的激活值 $a^{(2)}$ ）只使用一层计算单元（隐藏层）来得到

的。

在本节中，我们开始讨论深度神经网络，即含有多个隐藏层的神经网络。通过引入深度网络，我们可以计算更多复杂的输入特征。因为每一个隐藏层可以对上一层的输出进行非线性变换，因此深度神经网络拥有比“浅层”网络更加优异的表达能力（例如可以学习到更加复杂的函数关系）。

值得注意的是当训练深度网络的时候，每一层隐层应该使用非线性的激活函数 $f(x)$ 。这是因为多层的线性函数组合在一起本质上也只有线性函数的表达能力（例如，将多个线性方程组合在一起仅仅产生另一个线性方程）。因此，在激活函数是线性的情况下，相比于单隐藏层神经网络，包含多隐藏层的深度网络并没有增加表达能力。

深度网络的优势

为什么我们要使用深度网络呢？使用深度网络最主要的优势在于，它能以更加紧凑简洁的方式来表达比浅层网络大得多的函数集合。正式点说，我们可以找到一些函数，这些函数可以用 k 层网络简洁地表达出来（这里的简洁是指隐层单元的数目只需与输入单元数目呈多项式关系）。但是对于一个只有 $k - 1$ 层的网络而言，除非它使用与输入单元数目呈指数关系的隐层单元数目，否则不能简洁表达这些函数。

举一个简单的例子，比如我们打算构建一个布尔网络来计算 n 个输入比特的奇偶校验码（或者进行异或运算）。假设网络中的每一个节点都可以进行逻辑“或”运算（或者“与非”运算），亦或者逻辑“与”运算。如果我们拥有一个仅仅由一个输入层、一个隐层以及一个输出层构成的网络，那么该奇偶校验函数所需要的节点数目与输入层的规模 n 呈指数关系。但是，如果我们构建一个更深点的网络，那么这个网络的规模就可做到仅仅是 n 的多项式函数。

当处理对象是图像时，我们能够使用深度网络学习到“部分-整体”的分解关系。例如，第一层可以学习如何将图像中的像素组合在一起来检测边缘（正如我们在前面的练习中做的那样）。第二层可以将边缘组合起来检测更长的轮廓或者简单的“目标的部件”。在更深的层次上，可以将这些轮廓进一步组合起来以检测更为复杂的特征。

最后要提的一点是，大脑皮层同样是分多层进行计算的。例如视觉图像在人脑中是分多个阶段进行处理的，首先是进入大脑皮层的“V1”区，然后紧跟着进入大脑皮层“V2”区，以此类推。

训练深度网络的困难

虽然几十年前人们就发现了深度网络在理论上的简洁性和较强的表达能力，但是直到最近，研究者们也没有在训练深度网络方面取得多少进步。问题原因在于研究者们主要使用的学习算法是：首先随机初始化深度网络的权重，然后使用有监督的目标函数在有标签的训练集 $\left\{ \left(x_l^{(1)}, y^{(1)} \right), \dots, \left(x_l^{(m_l)}, y^{(m_l)} \right) \right\}$ 上进行训练。例如通过使用梯度下降法来降低训练误差。然而，这种方法通常不是十分奏效。这其中有如以下几方面原因：

数据获取问题

使用上面提到的方法，我们需要依赖于有标签的数据才能进行训练。然而有标签的数据通常是稀缺的，因此对于许多问题，我们很难获得足够多的样本来拟合一个复杂模型的参数。例如，考虑到深度网络具有强大的表达能力，在不充足的数据上进行训练将会导致过拟合。

局部极值问题

使用监督学习方法来对浅层网络（只有一个隐藏层）进行训练通常能够使参数收敛到合理的范围内。但是当用这种方法来训练深度网络的时候，并不能取得很好的效果。特别的，使用监督学习方法

训练神经网络时，通常会涉及到求解一个高度非凸的优化问题（例如最小化训练误差 $\sum_i \|h_W(x^{(i)}) - y^{(i)}\|^2$ ，其中参数 W 是要优化的参数。对深度网络而言，这种非凸优化问题的搜索区域中充斥着大量“坏”的局部极值，因而使用梯度下降法（或者像共轭梯度下降法，L-BFGS等方法）效果并不好。梯度弥散问题 梯度下降法（以及相关的L-BFGS算法等）在使用随机初始化权重的深度网络上效果不好的技术原因是：梯度会变得非常小。具体而言，当使用反向传播方法计算导数的时候，随着网络的深度的增加，反向传播的梯度（从输出层到网络的最初几层）的幅度值会急剧地减小。结果就造成了整体的损失函数相对于最初几层的权重的导数非常小。这样，当使用梯度下降法的时候，最初几层的权重变化非常缓慢，以至于它们不能够从样本中进行有效的学习。这种问题通常被称为“梯度的弥散”。

与梯度弥散问题紧密相关的问题是：当神经网络中的最后几层含有足够数量神经元的时候，可能单独这几层就足以对有标签数据进行建模，而不用最初几层的帮助。因此，对所有层都使用随机初始化的方法训练得到的整个网络的性能将会与训练得到的浅层网络（仅由深度网络的最后几层组成的浅层网络）的性能相似。

逐层贪婪训练方法

那么，我们应该如何训练深度网络呢？逐层贪婪训练方法是取得一定成功的一种方法。我们会在后面的章节中详细阐述这种方法的细节。简单来说，逐层贪婪算法的主要思路是每次只训练网络中的一层，即我们首先训练一个只含一个隐藏层的网络，仅当这层网络训练结束之后才开始训练一个有两个隐藏层的网络，以此类推。在每一步中，我们把已经训练好的前 $k - 1$ 层固定，然后增加第 k 层（也就是将我们已经训练好的前 $k - 1$ 的输出作为输入）。每一层的训练可以是有监督的（例如，将每一步的分类误差作为目标函

数），但更通常使用无监督方法（例如自动编码器，我们会在后边的章节中给出细节）。这些各层单独训练所得到的权重被用来初始化最终（或者说全部）的深度网络的权重，然后对整个网络进行“微调”（即把所有层放在一起优化有标签训练集上的训练误差）。逐层贪婪的训练方法取得成功要归功于以下几方面：

数据获取

虽然获取有标签数据的代价是昂贵的，但获取大量的无标签数据是容易的。自学习方法（self-taught learning）的潜力在于它能通过使用大量的无标签数据来学习到更好的模型。具体而言，该方法使用无标签数据来学习得到所有层（不包括用于预测标签的最终分类层） $W^{(l)}$ 的最佳初始权重。相比纯监督学习方法，这种自学习方法能够利用多得多的数据，并且能够学习和发现数据中存在的模式。因此该方法通常能够提高分类器的性能。

更好的局部极值 当用无标签数据训练完网络后，相比于随机初始化而言，各层初始权重会位于参数空间中较好的位置上。然后我们可以从这些位置出发进一步微调权重。从经验上来说，以这些位置为起点开始梯度下降更有可能收敛到比较好的局部极值点，这是因为无标签数据已经提供了大量输入数据中包含的模式的先验信息。

在下一节中，我们将会具体阐述如何进行逐层贪婪训练。

中英文对照

深度网络 Deep Networks

深度神经网络 deep neural networks

非线性变换 non-linear transformation

激活函数 activation function

简洁地表达 represent compactly

“部分-整体”的分解 part-whole decompositions

目标的部件 parts of objects

高度非凸的优化问题 highly non-convex optimization problem

共轭梯度 conjugate gradient

梯度的弥散 diffusion of gradients

逐层贪婪训练方法 Greedy layer-wise training

自动编码器 autoencoder

微调 fine-tuned

自学习方法 self-taught learning

6.3 栈式自编码算法

概述

逐层贪婪训练法依次训练网络的每一层，进而预训练整个深度神经网络。在本节中，我们将会学习如何将自编码器“栈化”到逐层贪婪训练法中，从而预训练（或者说初始化）深度神经网络的权重。

栈式自编码神经网络是一个由多层稀疏自编码器组成的神经网络，其前一层自编码器的输出作为其下一层自编码器的输入。对于一个 n 层栈式自编码神经网络，我们沿用自编码器一章的各种符号，假定用 $W^{(k,1)}, W^{(k,2)}, b^{(k,1)}, b^{(k,2)}$ 表示第 k 个自编码器对应的 $W^{(1)}, W^{(2)}, b^{(1)}, b^{(2)}$ 参数，那么该栈式自编码神经网络的编码过程就是，按照从前向后的顺序执行每一层自编码器的编码步骤：

$$a^{(l)} = f(z^{(l)}) \quad (6.1)$$

$$z^{(l+1)} = W^{(l,1)} a^{(l)} + b^{(l,1)} \quad (6.2)$$

同理，栈式神经网络的解码过程就是，按照从后向前的顺序执行每一层自编码器的解码步骤：

$$a^{(n+l)} = f(z^{(n+l)}) \quad (6.3)$$

$$z^{(n+l+1)} = W^{(n-l,2)} a^{(n+l)} + b^{(n-l,2)} \quad (6.4)$$

其中， $a^{(n)}$ 是最深层隐藏单元的激活值，其包含了我们感兴趣的信息，这个向量也是对输入值的更高阶的表示。

通过将 $a^{(n)}$ 作为softmax分类器的输入特征，可以将栈式自编码神经网络中学到的特征用于分类问题。

训练

一种比较好的获取栈式自编码神经网络参数的方法是采用逐层贪婪训练法进行训练。即先利用原始输入来训练网络的第一层，得到其参数 $W^{(1,1)}, W^{(1,2)}, b^{(1,1)}, b^{(1,2)}$ ；然后网络第一层将原始输入转化成为由隐藏单元激活值组成的向量（假设该向量为A），接着把A作为第二层的输入，继续训练得到第二层的参数 $W^{(2,1)}, W^{(2,2)}, b^{(2,1)}, b^{(2,2)}$ ；最后，对后面的各层同样采用的策略，即将前层的输出作为下一层输入的方式依次训练。

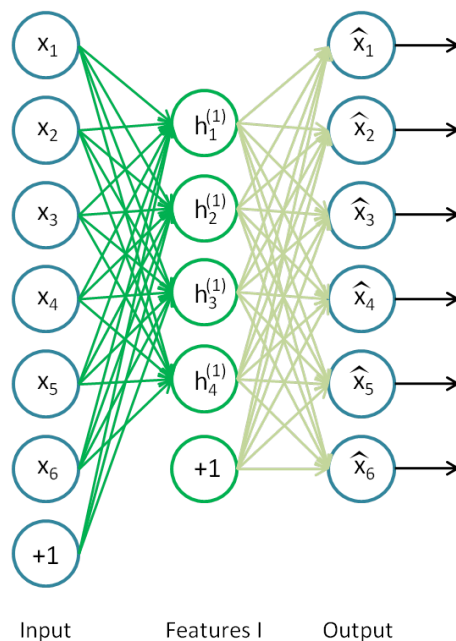
对于上述训练方式，在训练每一层参数的时候，会固定其它各层参数保持不变。所以，如果想得到更好的结果，在上述预训练过程完成之后，可以通过反向传播算法同时调整所有层的参数以改善结果，这个过程一般被称作“微调（fine-tuning）”。

实际上，使用逐层贪婪训练方法将参数训练到快要收敛时，应该使用微调。反之，如果直接在随机化的初始权重上使用微调，那么会得到不好的结果，因为参数会收敛到局部最优。

如果你只对以分类为目的的微调感兴趣，那么惯用的做法是丢掉栈式自编码网络的“解码”层，直接把最后一个隐藏层的 $a^{(n)}$ 作为特征输入到softmax分类器进行分类，这样，分类器（softmax）的分类错误的梯度值就可以直接反向传播给编码层了。

具体实例

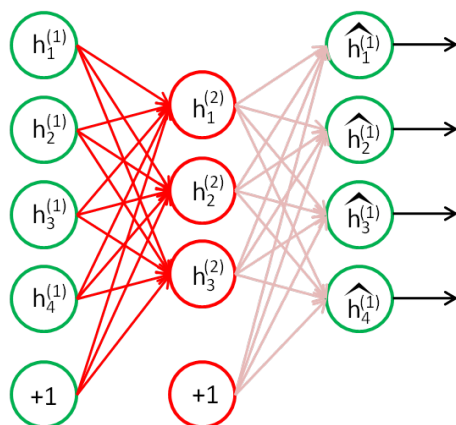
让我们来看个具体的例子，假设你想要训练一个包含两个隐含层的栈式自编码网络，用来进行MNIST手写数字分类（这将会是你的下一个练习）。首先，你需要用原始输入 $x^{(k)}$ 训练第一个自编码器，它能够学习得到原始输入的一阶特征表示 $h^{(1)(k)}$ （如下图所示）。



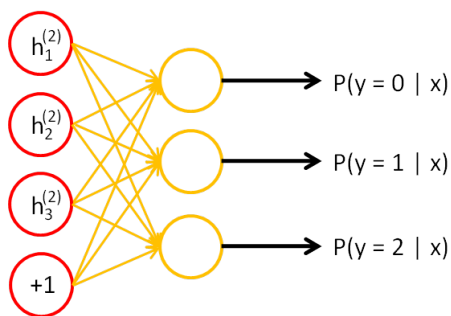
接着，你需要把原始数据输入到上述训练好的稀疏自编码器中，对于每一个输入 $x^{(k)}$ ，都可以得到它对应的一阶特征表示 $h^{(1)(k)}$ 。然后你再用这些一阶特征作为另一个稀疏自编码器的输入，使用它们来学习二阶特征 $h^{(2)(k)}$ 。（如下图所示）

同样，再把一阶特征输入到刚训练好的第二层稀疏自编码器中，得到每个 $h^{(1)(k)}$ 对应的二阶特征激活值 $h^{(2)(k)}$ 。接下来，你可以把这些二阶特征作为softmax分类器的输入，训练得到一个能将二阶特征映射到数字标签的模型。

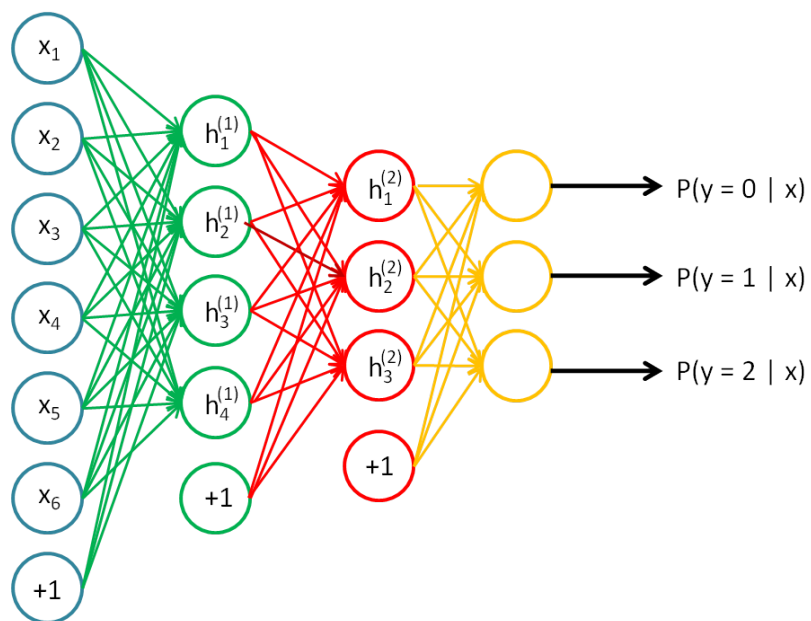
如下图所示，最终，你可以将这三层结合起来构建一个包含两个隐藏层和一个最终softmax分类器层的栈式自编码网络，这个网络能够如你所愿地对MNIST数字进行分类。



Input (Features I) Features II Output



Input (Features II) Softmax classifier



Input Features I Features II Softmax classifier

讨论

栈式自编码神经网络具有强大的表达能力及深度神经网络的所有优点。更进一步，它通常能够获取到输入的“层次型分组”或者“部分-整体分解”结构。为了弄清这一点，回顾一下，自编码器倾向于学习得到能更好地表示输入数据的特征。因此，栈式自编码神经网络的第一层会学习得到原始输入的一阶特征（比如图片里的边缘），第二层会学习得到二阶特征，该特征对应一阶特征里包含的一些模式（比如在构成轮廓或者角点时，什么样的边缘会共现）。栈式自编码神经网络的更高层还会学到更高阶的特征。

举个例子，如果网络的输入数据是图像，网络的第一层会学习如何去识别边，第二层一般会学习如何去组合边，从而构成轮廓、角等。更高层会学习如何去组合更形象且有意义的特征。例如，如果输入数据集包含人脸图像，更高层会学习如何识别或组合眼睛、鼻子、嘴等人脸器官。

中英文对照

自编码器 Autoencoder

逐层贪婪训练法 Greedy layer-wise training

预训练 PreTrain

栈式自编码神经网络 Stacked autoencoder

微调 Fine-tuning

原始输入 Raw inputs

层次型分组 Hierarchical grouping

部分-整体分解 Part-whole decomposition

一阶特征 First-order features

二阶特征 Second-order features

更高阶特征 Higher-order features

激活值 Activation

6.4 微调多层自编码算法

介绍

微调是深度学习中的常用策略，可以大幅提升一个栈式自编码神经网络的性能表现。从更高的视角来讲，微调将栈式自编码神经网络的所有层视为一个模型，这样在每次迭代中，网络中所有的权重值都可以被优化。一般策略

幸运的是，实施微调栈式自编码神经网络所需的工具都已齐备！为了在每次迭代中计算所有层的梯度，我们需要使用稀疏自动编码一节中讨论的反向传播算法。因为反向传播算法可以延伸应用到任意多层，所以事实上，该算法对任意多层的栈式自编码神经网络都适用。使用反向传播法进行微调

为方便读者，下面我们简要描述如何实施反向传播算法：

1. 进行一次前馈传递，对 L_2 层、 L_3 层直到输出层 L_{n_l} ，使用前向传播步骤中定义的公式计算各层上的激活值（激励响应）。

2. 对输出层（ n_l 层），令

$$\delta^{(n_l)} = -(\nabla_{a^{n_l}} J) \bullet f'(z^{(n_l)}) \quad (6.5)$$

（当使用softmax分类器时，softmax层满足： $\nabla J = \theta^T(I - P)$ ，其中 I 为输入数据对应的类别标签， P 为条件概率向量。）

3. 对 $l = n_l - 1, n_l - 2, n_l - 3, \dots, 2$ 令

$$\delta^{(l)} = \left((W^{(l+1)})^T \delta^{(l+1)} \right) \bullet f'(z^{(l)}) \quad (6.6)$$

4. 计算所需的偏导数：

$$\nabla_{W^{(l)}} J(W, b; x, y) = \delta^{(l+1)} (a^{(l)})^T, \quad (6.7)$$

$$\nabla_{b^{(l)}} J(W, b; x, y) = \delta^{(l+1)}. \quad (6.8)$$

$$J(W, b) = \left[\frac{1}{m} \sum_{i=1}^m J(W, b; x^{(i)}, y^{(i)}) \right] \quad (6.9)$$

注：我们可以认为输出层softmax分类器是附加上的一层，但是其求导过程需要单独处理。具体地说，网络“最后一层”的特征会进入softmax分类器。所以，第二步中的导数由 $\delta^{(n_l)} = -(\nabla_{a^{n_l}} J) \bullet f'(z^{(n_l)})$ 计算，其中 $\nabla J = \theta^T(I - P)$ 。

中英文对照

栈式自编码神经网络（可以考虑翻译为“多层自动编码器”或“多层自动编码神经网络”） Stacked autoencoder

微调 Fine tuning

反向传播算法 Backpropagation Algorithm

前馈传递 feedforward pass

激活值（可以考虑翻译为“激励响应”或“响应”） activation

Chapter 7

自编码线性解码器

7.1 线性解码器

稀疏自编码重述

稀疏自编码器包含3层神经元，分别是输入层，隐含层以及输出层。从前面（神经网络）自编码器描述可知，位于神经网络中的神经元都采用相同的激励函数。在注解中，我们修改了自编码器定义，使得某些神经元采用不同的激励函数。这样得到的模型更容易应用，而且模型对参数的变化也更为鲁棒。

回想一下，输出层神经元计算公式如下：

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)} \quad (7.1)$$

$$a^{(3)} = f(z^{(3)}) \quad (7.2)$$

其中 $a^{(3)}$ 是输出。在自编码器中， $a^{(3)}$ 近似重构了输入 $x = a^{(1)}$ 。

S 型激励函数输出范围是 $[0, 1]$ ，当 $f(z^{(3)})$ 采用该激励函数时，就要对输入限制或缩放，使其位于 $[0, 1]$ 范围中。一些数据集，比如 MNIST，能方便将输出缩放到 $[0, 1]$ 中，但是很难满足对输入值的要求。比如，PCA 白化处理的输入并不满足 $[0, 1]$ 范围要求，也不清楚是否有最好的办法可以将数据缩放到特定范围中。

线性解码器

设定 $a(3) = z(3)$ 可以很简单的解决上述问题。从形式上来看，就是输出端使用恒等函数 $f(z) = z$ 作为激励函数，于是有 $a(3) = f(z(3)) = z(3)$ 。我们称该特殊的激励函数为 **线性激励函数**（称为恒等激励函数可能更好些）。需要注意，神经网络中隐含层的神经元依然使用S型（或者tanh）激励函数。这样隐含单元的激励公式为 $a^{(2)} = \sigma(W^{(1)}x + b^{(1)})$ ，其中 $\sigma(\cdot)$ 是 S 型函数， x 是输入， $W(1)$ 和 $b(1)$ 分别是隐单元的权重和偏差项。我们仅在输出层中使用线性激励函数。一个 S 型或 tanh 隐含层以及线性输出层构成的自编码器，我们称为**线性解码器**。在这个线性解码器模型中， $\hat{x} = a^{(3)} = z^{(3)} = W^{(2)}a + b^{(2)}$ 。因为输出 \hat{x} 是隐单元激励输出的线性函数，改变 $W(2)$ ，可以使输出值 $a(3)$ 大于 1 或者小于 0。这使得我们可以用实值输入来训练稀疏自编码器，避免预先缩放样本到给定范围。随着输出单元的激励函数的改变，这个输出单元梯度也相应变化。回顾之前每一个输出单元误差项定义为：

$$\delta_i^{(3)} = \frac{\partial}{\partial z_i} \frac{1}{2} \|y - \hat{x}\|^2 = -(y_i - \hat{x}_i) \cdot f'(z_i^{(3)}) \quad (7.3)$$

其中 $y = x$ 是所期望的输出， \hat{x} 是自编码器的输出， $f(\cdot)$ 是激励函数。因为在输出层激励函数为 $f(z) = z$ ，这样 $f'(z) = 1$ ，所以上述公式可以简化为

$$\delta_i^{(3)} = -(y_i - \hat{x}_i) \quad (7.4)$$

当然，若使用反向传播算法来计算隐含层的误差项时：

$$\delta^{(2)} = \left((W^{(2)})^T \delta^{(3)} \right) \bullet f'(z^{(2)}) \quad (7.5)$$

因为隐含层采用一个 S 型（或 tanh）的激励函数 f ，在上述公式中， $f'(\cdot)$ 依然是 S 型（或 tanh）函数的导数。

中英文对照

线性解码器 Linear Decoders

稀疏自编码 Sparse Autoencoder

输入层 input layer

隐含层 hidden layer

输出层 output layer

神经元 neuron

神经网络 neural network

自编码器 autoencoder

激励函数 activation function

鲁棒 robust

S型激励函数 sigmoid activation function

tanh激励函数 tanh function

线性激励函数 linear activation function

恒等激励函数 identity activation function

隐单元 hidden unit

权重 weight

偏差项 error term

反向传播算法 backpropagation

Chapter 8

处理大型图像

8.1 卷积特征提取

概述

前面的练习中，解决了一些有关低分辨率图像的问题，比如：小块图像，手写数字小幅图像等。在这部分中，我们将把已知的方法扩展到实际应用中更加常见的大图像数据集。

全联通网络

在稀疏自编码章节中，我们介绍了把输入层和隐含层进行“全连接”的设计。从计算的角度来讲，在其他章节中曾经用过的相对较小的图像（如在稀疏自编码的作业中用到过的 8×8 的小块图像，在MNIST数据集中用到过的 28×28 的小块图像），从整幅图像中计算特征是可行的。但是，如果是更大的图像（如 96×96 的图像），要通过这种全联通网络的这种方法来学习整幅图像上的特征，从计算角度而言，将变得非常耗时。你需要设计 10^4 （=10000）个输入单元，假设你要学习 100 个特征，那么就有 10^6 个参数需要去学习。与 28×28 的小块图像相比较， 96×96 的图像使用前向输送或者后向传导的计算方式，计算过程也会慢 10^2 （=100）倍。

部分联通网络

解决这类问题的一种简单方法是对隐含单元和输入单元间的连接

加以限制：每个隐含单元仅仅只能连接输入单元的一部分。例如，每个隐含单元仅仅连接输入图像的一小片相邻区域。（对于不同于图像输入的输入形式，也会有一些特别的连接到单隐含层的输入信号“连接区域”选择方式。如音频作为一种信号输入方式，一个隐含单元所需要连接的输入单元的子集，可能仅仅是一段音频输入所对应的某个时间段上的信号。）

网络部分连通的思想，也是受启发于生物学里面的视觉系统结构。视觉皮层的神经元就是局部接受信息的（即这些神经元只响应某些特定区域的刺激）。

卷积

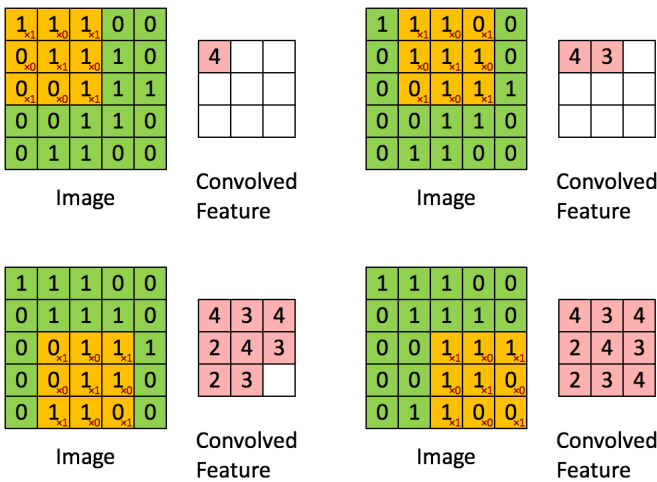
自然图像有其固有特性，也就是说，图像的一部分的统计特性与其他部分是一样的。这也意味着我们在这一部分学习的特征也能用在另一部分上，所以对于这个图像上的所有位置，我们都能使用同样的学习特征。

更恰当的解释是，当从一个大尺寸图像中随机选取一小块，比如说 8×8 作为样本，并且从这个小块样本中学习到了某些特征，这时我们可以把从这个 8×8 样本中学习到的特征作为探测器，应用到这个图像的任意地方中去。特别是，我们可以用从 8×8 样本中所学习到的特征跟原本的大尺寸图像作卷积，从而对这个大尺寸图像上的任一位置获得一个不同特征的激活值。

下面给出一个具体的例子：假设你已经从一个 96×96 的图像中学习到了它的一个 8×8 的样本所具有的特征，假设这是由有 100 个隐含单元的自编码完成的。为了得到卷积特征，需要对 96×96 的图像的每个 8×8 的小块图像区域都进行卷积运算。也就是说，抽取 8×8 的小块区域，并且从起始坐标开始依次标记为 $(1, 1)$ ， $(1, 2)$ ，...，一直到 $(89, 89)$ ，然后对抽取的区域逐个运行训练过的稀疏自编码来得到特征的激活值。在这个例子里，显然可以

得到 100 个集合，每个集合含有 89x89 个卷积特征。

假设给定了 $r \times c$ 的大尺寸图像，将其定义为 x_{large} 。首先通过从大尺寸图像中抽取的 $a \times b$ 的小尺寸图像样本 x_{small} 训练稀疏自编码，计算 $f = (W(1)x_{small} + b(1))$ (σ 是一个 sigmoid 型函数) 得到了 k 个特征，其中 $W(1)$ 和 $b(1)$ 是可视层单元和隐含单元之间的权重和偏差值。对于每一个 $a \times b$ 大小的小图像 x_s ，计算出对应的值 $f_s = (W(1)x_s + b(1))$ ，对这些 $f_{convolved}$ 值做卷积，就可以得到 $k \times (r - a + 1) \times (c - b + 1)$ 个卷积后的特征的矩阵。



在接下来的章节里，我们会更进一步描述如何把这些特征汇总到一起以得到一些更利于分类的特征。

中英文对照

- 全联通网络 Full Connected Networks
- 稀疏编码 Sparse Autoencoder
- 前向输送 Feedforward
- 反向传播 Backpropagation
- 部分联通网络 Locally Connected Networks
- 连接区域 Contiguous Groups
- 视觉皮层 Visual Cortex
- 卷积 Convolution

固有特征 Stationary 池化 Pool

8.2 池化

概述

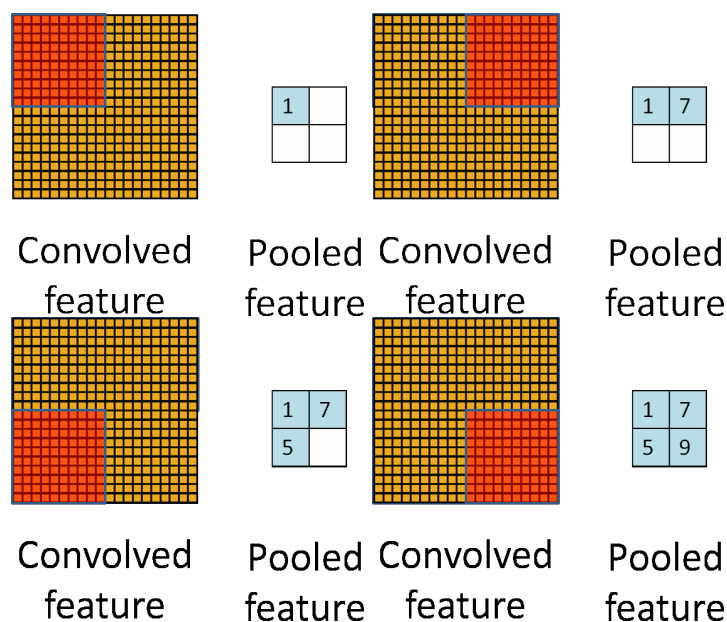
在通过卷积获得了特征 (features) 之后, 下一步我们希望利用这些特征去做分类。理论上讲, 人们可以用所有提取得到的特征去训练分类器, 例如 softmax 分类器, 但这样做面临计算量的挑战。例如: 对于一个 96×96 像素的图像, 假设我们已经学习得到了 400 个定义在 8×8 输入上的特征, 每一个特征和图像卷积都会得到一个 $(96 - 8 + 1) * (96 - 8 + 1) = 7921$ 维的卷积特征, 由于有 400 个特征, 所以每个样例 (example) 都会得到一个 $7921 * 400 = 3,168,400$ 维的卷积特征向量。学习一个拥有超过 3 百万特征输入的分类器十分不便, 并且容易出现过拟合 (over-fitting)。

为了解决这个问题, 首先回忆一下, 我们之所以决定使用卷积后的特征是因为图像具有一种“静态性”的属性, 这也就意味着在一个图像区域有用的特征极有可能在另一个区域同样适用。因此, 为了描述大的图像, 一个很自然的想法就是对不同位置的特征进行聚合统计, 例如, 人们可以计算图像一个区域上的某个特定特征的平均值 (或最大值)。这些概要统计特征不仅具有低得多的维度 (相比使用所有提取得到的特征), 同时还会改善结果 (不容易过拟合)。这种聚合的操作就叫做池化 (pooling), 有时也称为平均池化或者最大池化 (取决于计算池化的方法)。

下图显示池化如何应用于一个图像的四块不重合区域。

池化的不变性

如果人们选择图像中的连续范围作为池化区域, 并且只是池化相同 (重复) 的隐藏单元产生的特征, 那么, 这些池化单元就具有平移不变性 (translation invariant)。这就意味着即使图像经历了一个小的



平移之后，依然会产生相同的 (池化的) 特征。在很多任务中 (例如物体检测、声音识别)，我们都更希望得到具有平移不变性的特征，因为即使图像经过了平移，样例(图像)的标记仍然保持不变。例如，如果你处理一个MNIST数据集的数字，把它向左侧或右侧平移，那么不论最终的位置在哪里，你都会期望你的分类器仍然能够精确地将其分类为相同的数字。 (*MNIST 是一个手写数字库识别库: <http://yann.lecun.com/exdb/mnist/>)

形式化描述

形式上，在获取到我们前面讨论过的卷积特征后，我们要确定池化区域的大小(假定为 $m \times n$)，来池化我们的卷积特征。那么，我们把卷积特征划分到数个大小为 $m \times n$ 的不相交区域上，然后用这些区域的平均(或最大)特征来获取池化后的卷积特征。这些池化后的特征便可以用来做分类。

中英文对照

特征 features

样例 example

过拟合 over-fitting

平移不变性 translation invariant

池化 pooling

提取 extract

物体检测 object detection