

深度学习算法掘金 ALPHA 因子

大数据深度学习系列之二

报告摘要:

● 金融大数据下的 Alpha 因子挖掘

多因子 Alpha 策略是发掘出驱动个股产生 Alpha 收益的因子, 根据有效的 Alpha 因子设计相应的选股策略, 筛选投资的股票组合, 以寻找超越市场的股票超额收益。为了获取新的 Alpha 来源, 我们一方面可以对传统因子进行更加深入的挖掘, 例如挖掘因子的非线性特征, 寻找有效的因子组合。另一方面, 我们可以利用更加高效的数据挖掘手段从市场数据中间寻找新的 Alpha 因子。随着大数据时代的来临, 数据挖掘的方法不断革新改进, 浩如烟海的市场数据为 Alpha 的来源提供了巨大的可能性。

● 深度学习股价预测模型

从市场微观结构的角度来说, 股票价格的形成和变化是由买卖双方的交易行为决定的, 因此, 对高频市场行情数据的挖掘有可能获得对未来股票价格走势的有预测能力的模式。

本报告通过样本内大量历史数据训练深度学习预测模型, 对以周为频率的中证 800 股票价格涨跌进行预测, 建立起了可以对股价短期内走势进行预测的机器学习模型。

● 深度学习股票多因子交易策略

基于深度学习股价预测模型对股票价格变化的预测得分, 本报告提出了股票交易的 Alpha 策略。在组合规模为 100 的情况下, 该多因子 Alpha 策略自 2011 年以来累积收益率超过 120%, 各年度收益率都超过 15%。

● 结论

通过中证 800 成份股的实证研究, 本报告验证了深度学习这一大数据时代的机器学习利器在股票价格预测上的有效性。通过深度学习模型对市场数据进行挖掘, 获得了可以产生超额收益的因子, 该因子的表现超越了传统的 Alpha 因子。

图 1 不同规模的深度学习因子组合

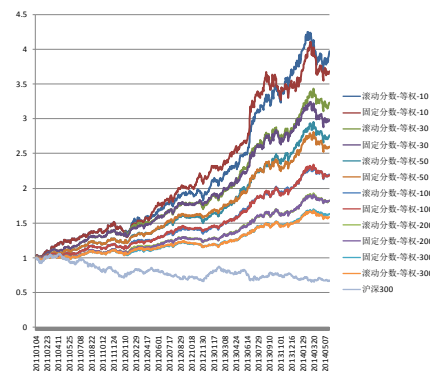
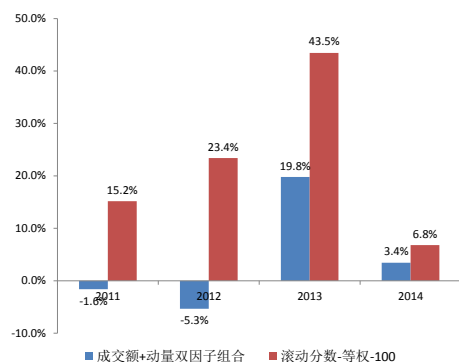


图 2 深度学习因子与价量双因子组合对比



分析师: 安宁宁 S0260512020003



0755-23948352



ann@gf.com.cn

相关研究:

大数据深度学习系列之一: 深 2014-06-18

深度学习之股指期货日内交易策略

目录索引

一、ALPHA 因子挖掘.....	4
（一）主流因子表现回顾.....	4
（二）大数据 ALPHA 因子挖掘.....	6
二、深度学习介绍.....	7
（一）深度学习：机器学习的新浪潮.....	7
（二）模型起源.....	9
（三）深层模型结构.....	10
（四）人工神经网络.....	11
（五）自编码器和深度网络.....	13
三、实证分析.....	16
（一）深度学习预测模型.....	17
（二）交易策略.....	18
（三）预测模型效果.....	19
（四）对冲策略效果.....	20
四、总结与讨论.....	26

图表索引

图 1: 基于多因子的 Alpha 策略框架	4
图 2: 主流 Alpha 因子自 2011 年 1 月至 2014 年 5 月的累积收益曲线.....	5
图 3: 主流 Alpha 因子的 IC (自 2011 年 1 月至 2014 年 5 月)	6
图 4: 机器学习的一般流程	8
图 5: 视觉系统的层级结构	10
图 6: 深度学习的层级结构	10
图 7: 神经元示意图	11
图 8: 逻辑函数输入输出图	12
图 9: 神经网络示意图	12
图 10: 深度学习示意图	14
图 11: 自编码器示意图	15
图 12: 降噪自编码器示意图	15
图 13: 基于深度学习的 Alpha 因子策略示意图	16
图 14: 基于深度学习的滚动预测模型示意图	17
图 15: 不同规模股票等权组合策略收益曲线	21
图 16: 深度学习对冲组合与成交金额-股价动量双因子对冲组合收益曲线对比....	22
图 17: 深度学习对冲组合与成交金额-股价动量双因子对冲组合月度收益率对比..	23
图 18: 深度学习不同预测因子及其组合的五档多空累计表现对比.....	24
图 19: 深度学习 Alpha 策略股票组合行业分布图	25
图 20: 深度学习 Alpha 策略样本内各期股票组合行业数量中位数.....	25
表 1: 深度学习在机器学习领域的重大突破	8
表 2: 深度学习股价预测特征选取	18
表 3: 深度学习股价预测结果	20
表 4: 不同数量股票等权组合策略年度收益率	22
表 5: 深度学习不同因子的表现对比	24

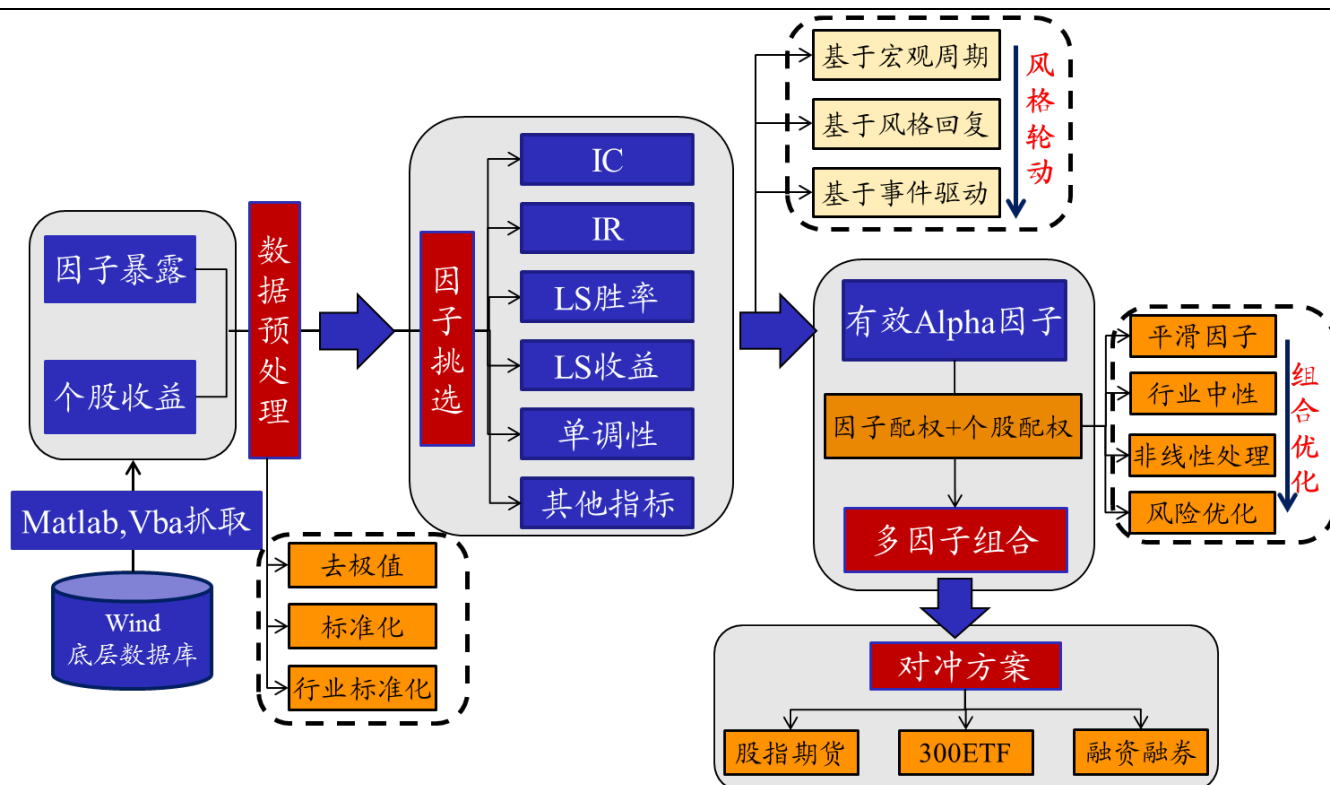
一、Alpha 因子挖掘

（一）主流因子表现回顾

现代金融理论则认为，证券投资者所获得的收益分为两部分：来自市场的平均收益（即 Beta 收益）以及独立于市场的超额收益（即 Alpha 收益）。近年来，随着大牛市行情的消失以及结构性行情的到来，个股出现严重分化的现象屡出不穷，许多投资者渐渐把目光投向了个股的 Alpha 收益，因此专门寻找 Alpha 收益的 Alpha 策略也开始兴起，相比传统的投资策略，阿尔法策略具有更强的主动性，投资者不再被动等待交易时机，而是通过主动选取具备 Alpha 正收益的股票，随时进场进行交易。由于 Alpha 策略是一种中性策略，因此投资者在弱势和振荡行情中也能获得稳定收益，不需要判断大盘走势。

多因子 Alpha 策略通过发掘出驱动个股产生 Alpha 收益的因子，根据有效的 Alpha 因子设计相应的选股策略，筛选投资的股票组合，以寻找超越市场的股票超额收益。基于多因子的 Alpha 策略框架如图 1 所示。注意到，在数据预处理阶段，为了消除不同行业不同股票本身的差异性，我们需要对因子进行标准化或者行业标准化。

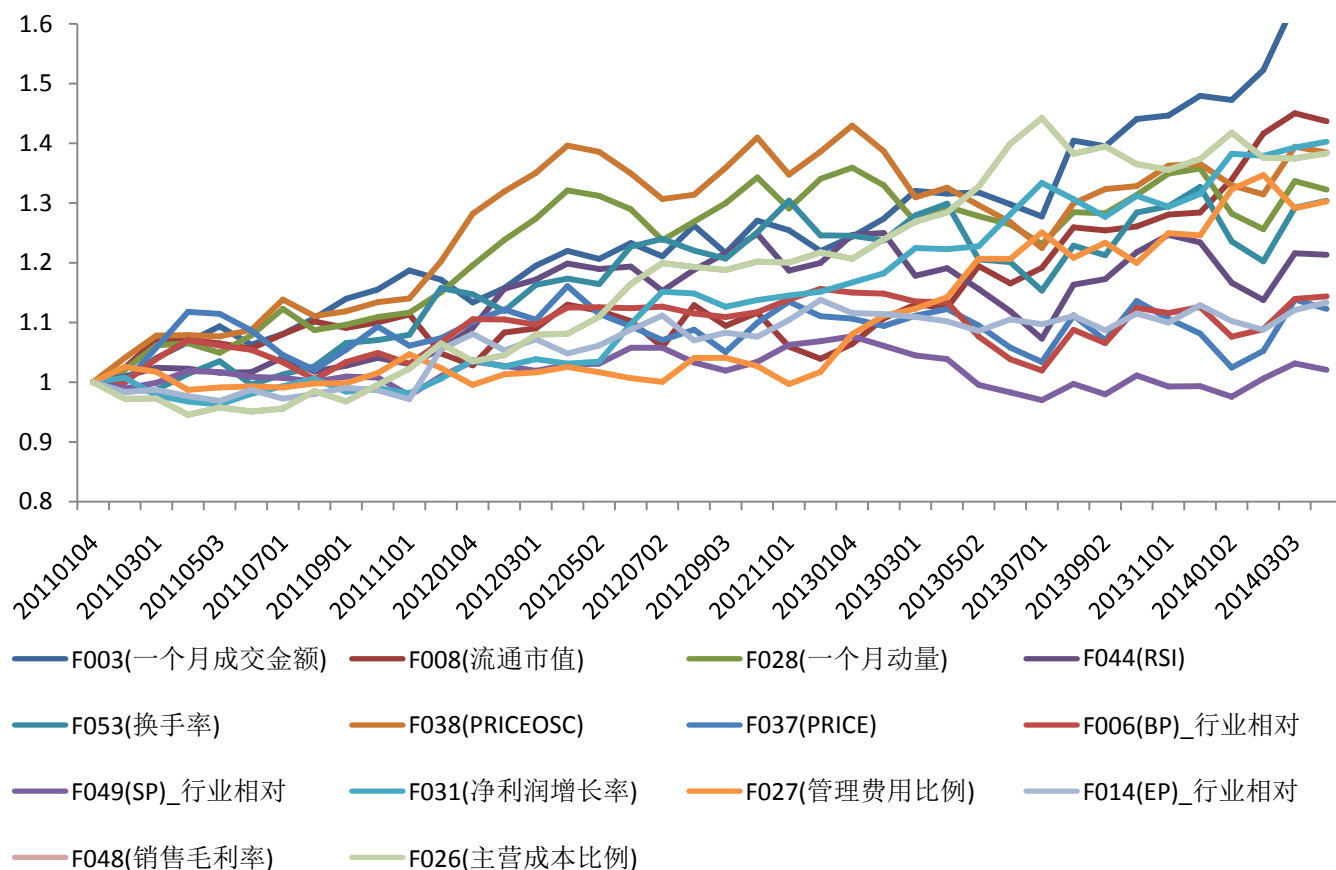
图1：基于多因子的Alpha策略框架



数据来源：广发证券发展研究中心

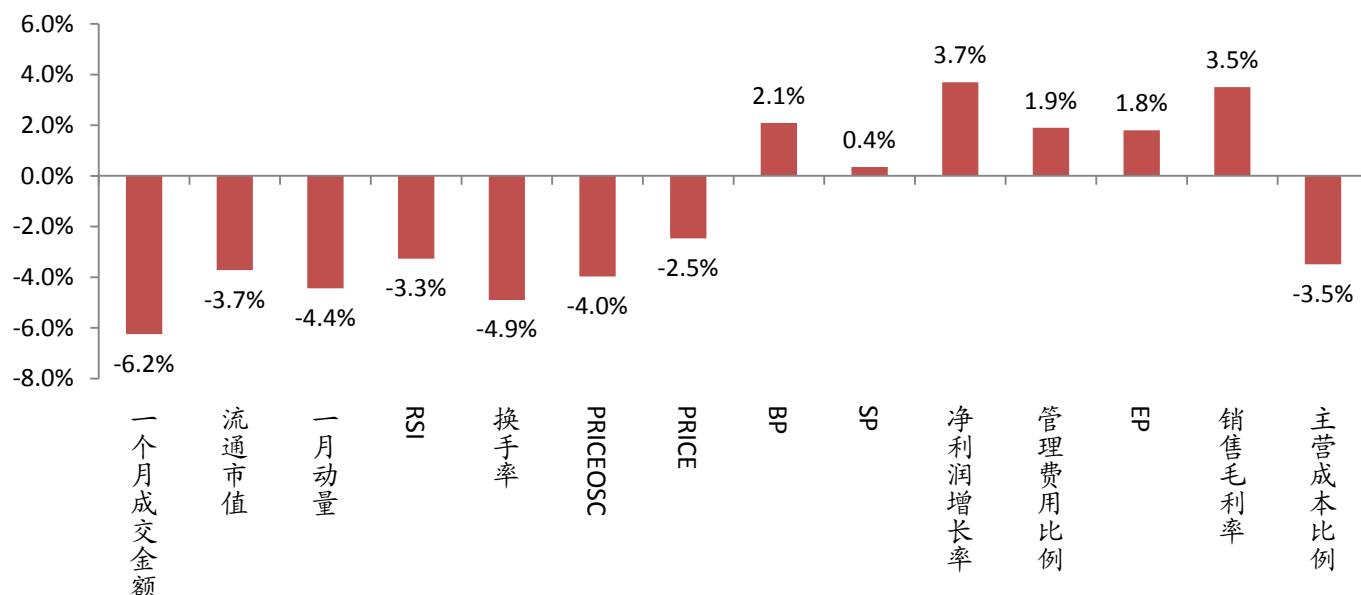
在中证 800 成份股的备选池下,考察自 11 年 1 月至 14 年 5 月主流 Alpha 因子的表现。可以看到如图 2 所示,在按月调仓下,规模因子、反转因子、成交量因子是最有效因子,股价、股价波动因子以及估值类因子次之。图 3 表示了各 Alpha 因子的 IC 的对比图。

图2: 主流Alpha因子自2011年1月至2014年5月的累积收益曲线



数据来源: 广发证券发展研究中心, Wind 数据库

图3：主流Alpha因子的IC（自2011年1月至2014年5月）



数据来源：广发证券发展研究中心，Wind 数据库

（二）大数据 Alpha 因子挖掘

为了获取新的 Alpha 来源，我们一方面可以对传统 Alpha 因子进行更加深入的挖掘，例如挖掘因子的非线性特征，或者是基于多个因子来寻找有效的因子组合。另一方面，我们可以利用更加高效的数据挖掘手段从市场数据中寻找新的 Alpha 因子。

从市场微观结构的角度来说，股票价格的形成和变化是由买卖双方的交易行为决定的，因此，对市场数据，特别是高频市场行情数据的挖掘有可能获得对未来股票价格走势的有预测能力的模式。

随着计算机科学与技术的蓬勃发展，存储成本的降低，计算速度的提高，人们越来越关注“云计算”、“大数据”这些热门词汇。金融市场的数据量也越来越多——单单 A 股市场就有近 2000 只股票，光考虑以秒为单位的高频数据，每个交易日就会产生两千多万新的数据样本。这为我们从市场数据获取 Alpha 来源提供了广阔的平台和无限的可能性。同时也带来了新的问题：如何从中“数据海洋”中提取有用的信息，来帮助投资者获得超额收益呢？这就需要借助于符合大数据时代要求的新一代机器学习算法了。

大数据时代，对机器学习和数据挖掘方法表达信息能力的要求也越来越高。传统的机器学习算法通常受囿于表达能力不强的问题，不能充分利用海量数据的信息。近年来在机器学习的科学研究和工业应用领域流行的深度学习算法有效的解决了这一问题，引

领了大数据时代机器学习的潮流。通过深度学习算法，能否从海量的市场数据中获得新的有效因子，在股票投资中获得超额收益呢？

本报告中，我们考虑基于高频的股票交易量价信息，建立起深度学习预测模型，对股票未来的涨跌情况进行预测。并且将预测结果作为因子，寄望于通过该因子获得 Alpha 收益。

二、深度学习介绍

（一）深度学习：机器学习的新浪潮

2012 年 6 月，《纽约时报》披露了“谷歌大脑”（Google Brain）项目，吸引了公众的广泛关注。这个项目是由著名的斯坦福大学机器学习教授吴恩达（Andrew Ng）和在大规模计算机系统方面的世界顶尖专家 Jeff Dean 共同主导，用 16000 个 CPU Core 的并行计算平台训练一种称为“深层神经网络”的机器学习模型，在语音识别和图像识别等领域获得了巨大的成功。在谷歌的一次公开展示中，该大脑从 1 千万个随机挑选的没有经过标注处理的 YouTube 视频中，自动识别出了猫脸。不久之后，谷歌收购了深度学习的教父级人物——加拿大多伦多大学教授 Geoffrey Hinton 创建的人工智能研究机构 DNNresearch，继续增大在深度学习领域的投入。

微软和 IBM 的研究人员使用深度学习在语音识别上也取得了巨大进展。2012 年 11 月，微软首席科学家 Richard Rashid 在中国天津的一次活动上公开演示了一个全自动的同声传译系统，讲演者用英文演讲，后台的计算机一气呵成自动完成语音识别、英中机器翻译，以及中文语音合成，效果非常流畅。后面支撑的关键技术就是深度学习。根据微软和谷歌的报告，用深度学习改进传统的隐马尔科夫语音识别模型，将语音识别的错误率相对降低了 30%。同时，深度学习技术在图像识别领域取得惊人的效果，2012 年在业界著名的 ImageNet 评测上将错误率从此前的最好成绩 26% 降低到 15%。也是在这一年，深度学习还被应用于制药公司的药物活性预测问题，并获得世界最好成绩。

另一家美国社交网络巨头 Facebook 也在 2013 年下半年组建了深度学习研究小组，用来分析预测用户的行为习惯。Facebook 此前已经使用传统的机器学习算法来给用户定向投递新闻和广告（谷歌，亚马逊，百度，阿里巴巴等互联网企业也正在使用这种针对具体用户的“推荐系统”），而今他们寄望于通过深度学习来获得更好的效果。事实上，在此前著名的 Netflix（奈飞公司，全球最大的流媒体播放服务提供商，《纸牌屋》的出品方）电影推荐系统比赛中，基于深度学习的算法就一举夺魁。

国内的互联网 IT 企业也开始了在深度学习领域的投入和研发。2013 年 1 月，在百度的年会上，创始人兼 CEO 李彦宏高调宣布要成立深度学习研究院，这是百度成立十多

年以来第一次成立研究院。同年，国内语音识别技术的领军企业科大讯飞也将深度学习作为未来研发的重点目标。

2013年4月，《麻省理工学院技术评论》（MIT Technology Review）杂志将深度学习列为2013年的十大突破性技术之首。深度学习在机器学习各领域的突破性成果如表1所示。

表1：深度学习在机器学习领域的重大突破

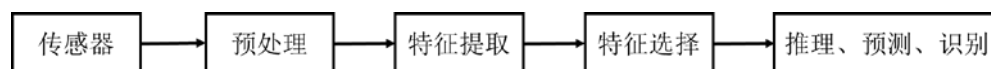
语音识别	识别错误率相对降低30%，是近20年最大突破
图像识别	ImageNet 评测上的错误率从以前方法的26%提升到15%
智能推荐系统	奈飞视频推荐比赛最佳成绩
药物活性预测	比赛最佳成绩
自然语言处理	准确率提高到85%，目前最好结果
CRT 网络点击率预估模型	目前最好结果

数据来源：广发证券发展研究中心

近年来深度学习相关的研究和应用如此火爆，那么，深度学习究竟是一种什么样的学习方式呢？

机器学习的目标有图像识别、语音识别、自然语言理解、股价预测、天气预测、基因表达、内容推荐等等。以图像识别为例，目前我们通过机器学习去解决这些问题的思路一般是如图4所示。

图4：机器学习的一般流程



数据来源：广发证券发展研究中心

从开始的通过传感器（例如CMOS）来获得数据。然后经过数据预处理、特征提取、特征选择，再到推理、预测或者识别。最后一个部分“推理、预测、识别”，也就是我们通常所说的机器学习的部分，主要是通过数学和统计方法来建立学习模型。而中间的三部分概括起来就是“特征表达”。良好的特征表达，对最终算法的准确性起了非常关键的作用，而整个机器学习系统建立时主要的计算和测试工作都耗在这一大部分。但实际中，这一块一般都是人工完成的，也就是靠人工提取特征。在股价预测中，类似的流程就是我们定义了很多的技术指标，即机器学习中的特征，比如移动平均线、布林线、相对强弱指标等；然后根据我们要解决的问题选择合适的特征来建立预测规则（模型）。

手工地选取特征是一件非常费力、启发式（利用专业知识来选取特征）的方法，不能选取好很大程度上靠经验和运气，而且它的调节需要大量的时间。既然手工选取特征不太好，那么能不能自动地学习一些特征呢？答案是肯定的！深度学习就是用来干这个事情的。

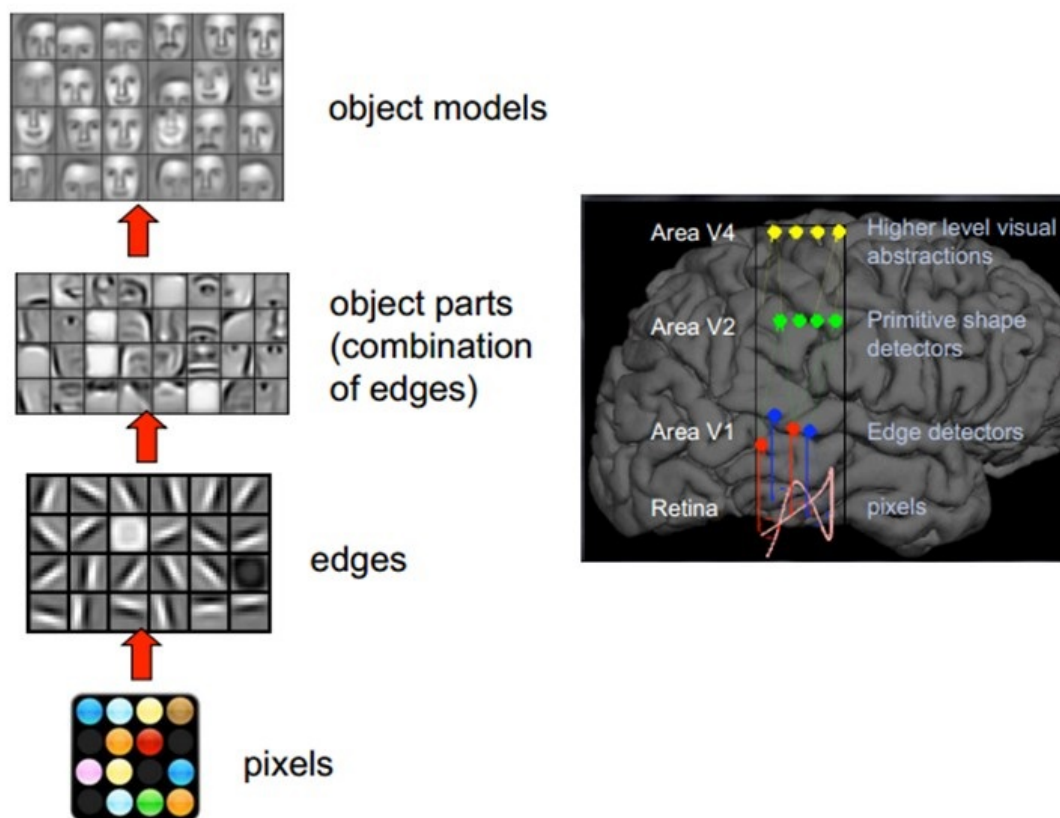
（二）模型起源

很多机器学习模型受到生物学方面的启发，比如说遗传算法，粒子群算法，蚁群算法等。深度学习也与生理医学上的发现有关。

1981 年的诺贝尔医学奖颁发给了 David Hubel 和 Torsten Wiesel，以及 Roger Sperry。前两位的主要贡献，是“发现了视觉系统的信息处理”，即 Hubel-Wiesel 模型。1958 年，Hubel 和 Wiesel 在约翰霍普金斯大学研究瞳孔区域与大脑皮层神经元的对应关系。他们在猫的后脑头骨上，开了一个 3 毫米的小洞，向洞里插入电极，测量神经元的活跃程度。然后，他们在小猫的眼前，展现各种形状、各种亮度的物体。并且，在展现每一件物体时，还改变物体放置的位置和角度。他们期望通过这个办法，让小猫瞳孔感受不同类型、不同强弱的刺激。之所以做这个试验，目的是去证明一个猜测：位于后脑皮层的不同视觉神经元，与瞳孔所受刺激之间，存在某种对应关系。一旦瞳孔受到某一种刺激，后脑皮层的某一部分神经元就会活跃。经历了很多天反复的枯燥的试验，Hubel 和 Wiesel 发现了一种被称为“方向选择性细胞”（Orientation Selective Cell）的神经元细胞。当瞳孔发现了眼前的物体的边缘，而且这个边缘指向某个方向时，这种神经元细胞就会活跃。

这个发现激发了人们对于神经系统的进一步思考：脑神经系统具有丰富的层次结构。神经-中枢-大脑的工作过程，或许是一个不断迭代、不断抽象的过程。这里的关键词有两个，一个是抽象，一个是迭代。从原始信号，做低级抽象，逐渐向高级抽象迭代。人类的逻辑思维，经常使用高度抽象的概念。如图 5 所示，从原始信号摄入开始（瞳孔摄入像素），接着做初步处理（V1 层：大脑皮层某些细胞发现边缘和方向），然后抽象（V2 层：大脑判定，眼前的物体的形状，是圆形的），然后进一步抽象（V3 层：大脑进一步判定该物体是人脸）。换句话说，人的视觉系统的信息处理是分级的。高层的特征是低层特征的组合，从低层到高层的特征表示越来越抽象，越来越能表现语义或者意图。而抽象层面越高，存在的可能猜测就越少，就越利于分类。例如，单词集合和句子的对应是多对一的，句子和语义的对应又是多对一的，语义和意图的对应还是多对一的，这是个层级体系。

图5：视觉系统的层级结构

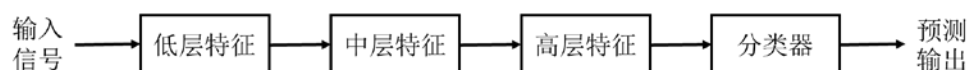


数据来源：广发证券发展研究中心

（三）深层模型结构

深度学习是模拟大脑皮层的 Hubel-Wiesel 模型，采用一层层“抽象化”的方式来对数据或者信号进行表达。类似于大脑皮层对图像的分辨，深度学习模型首先从原始信号（类似于人脸识别系统中的像素）中分离出低层的特征（类似人脸识别系统中物体的边），然后从低层特征中获取高一层的特征（类似于人脸识别系统中由边组成的轮廓），然后获得更高一层的表达（类似人脸识别中的人脸），最后在高层特征上建立起分类器，获得模型的预测输出。

图6：深度学习的层级结构



数据来源：广发证券发展研究中心

深度学习是在对大量的数据进行特征抽象的同时，获得其丰富的表达，而对于特定的学习目标，相应的，合适的表达会被激活，从而在没有经过人工特征选取的前提下获取足够好的学习效果。

深度学习中，先是采用逐层学习的贪婪式算法对特征进行提取，是一种无监督学习，即采用未经人工标注类别的样本进行学习（对应的，支持向量机，神经网络等方法有定义好的“输入”和“输出”，属于有监督学习）。逐层学习时，有多种选择方式，目前比较普遍的有自编码器和受限玻尔兹曼机，本报告中采用的模型是用自编码器进行学习的。这两种模型都是基于人工神经网络来实现的。因此在介绍具体的深度学习模型之前，我们先回顾一下人工神经网络的基本知识。

（四）人工神经网络

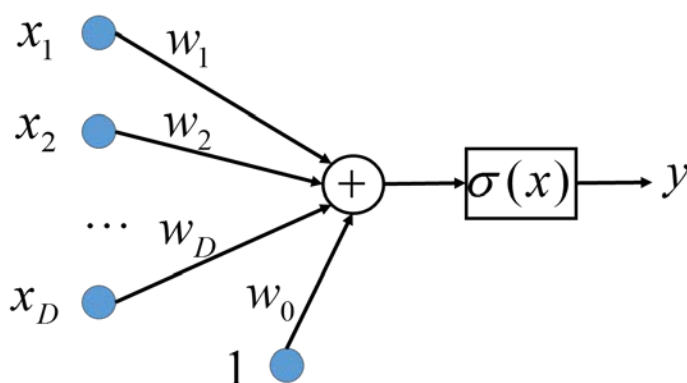
人工神经网络是一种应用类似于大脑神经突触连接的结构进行信息处理的数学模型，工程上常简称为神经网络。神经网络由大量的节点（或称“神经元”）和节点之间的相互连接构成。每个节点代表一种特定的输出函数，称为激励函数。每两个节点间的连接都代表一个对于通过该连接信号的加权值，称之为权重。网络的输出则依网络的连接方式、权重值和激励函数的不同而不同。

图 7 表示的是神经元的基本形式，将 D 个输入变量和偏置项累加起来，经过激励函数，获得输出 y 。常用的激励函数有逻辑函数，正切函数等。逻辑函数作为激励函数的表达式为

$$y = \sigma(x) = \frac{1}{1 + e^{-x}} \quad (1)$$

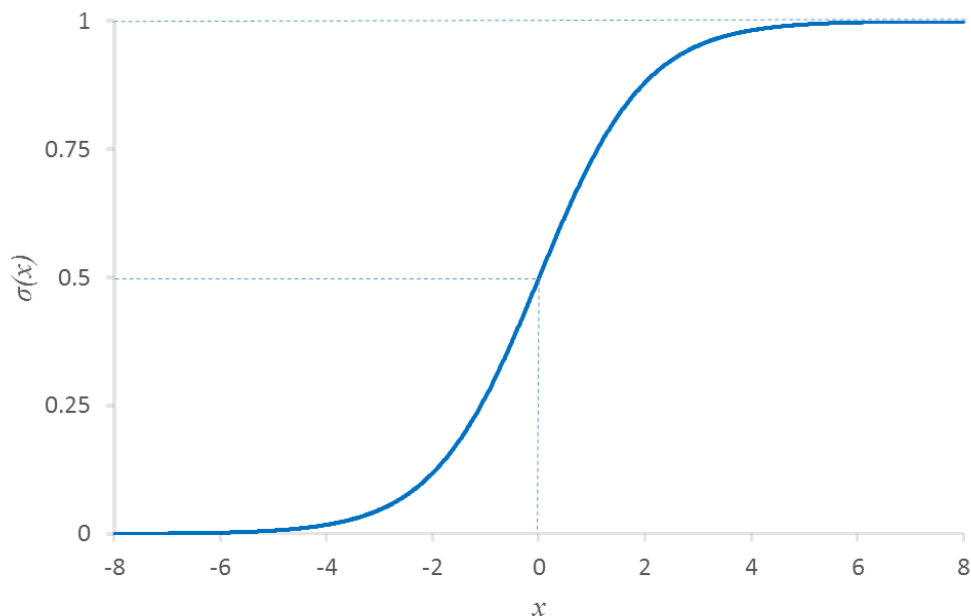
如图 8 所示。机器学习中常用的逻辑回归模型（Logistic Regression）就是采取这种形式的输出函数以达到分类的目的。

图 7：神经元示意图



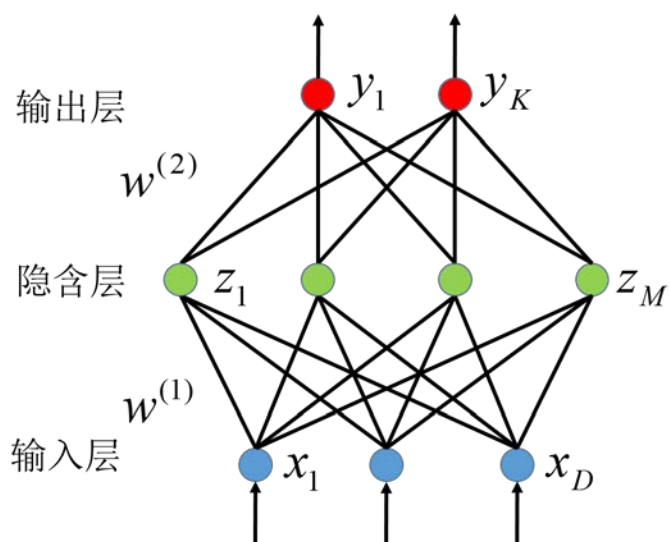
数据来源：广发证券发展研究中心

图8: 逻辑函数输入输出图



数据来源: 广发证券发展研究中心

图9: 神经网络示意图



数据来源: 广发证券发展研究中心

因此, 该神经元的数学表达式为

$$y = \sigma\left(\sum_{i=1}^D w_i x_i + w_0\right) \quad (2)$$

一个完整的神经网络模型通常将节点分成若干层次：输入层，输出层和隐含层，如图 9 所示。输入层即我们给定的模型输入，输出层即我们想通过神经网络“预测”的结果，隐含层相当于网络系统的状态。对于回归神经网络，输出层节点的个数即我们所要预测的变量个数；对于分类神经网络，输出层节点的个数通常是可能的分类总类别数。该神经网络第 k 个输出的数学表达式为

$$y_k = \sigma \left\{ \sum_{j=1}^M (w_{kj}^{(2)} h(\sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)}) + w_{k0}^{(2)}) \right\} \quad (3)$$

其中 σ 和 h 分别为输出层和隐含层的激励函数。神经网络的参数为各层的网络系数 w_{ij} ，可以一并记为向量 \mathbf{w} 。注：公式中的黑体表示向量或者矩阵， w_{ij} 为向量 \mathbf{w} 中的元素。

神经网络模型的学习即利用我们已经有的输入输出数据（训练集），对参数 \mathbf{w} 的优化，使得输出 y 尽可能的接近于其真实值 t ，即要使得如下的预测误差（即损失函数）最小化

$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) = \sum_{n=1}^N \sum_{k=1}^K (y_{nk} - t_{nk})^2 \quad (4)$$

一般使用梯度下降法来获取最优的参数 \mathbf{w} ：

$$w_{ij} := w_{ij} - \alpha \frac{\partial}{\partial w_{ij}} E(\mathbf{w}) \quad (5)$$

其中参数 α 为学习率，表示每一次迭代的步长。当神经网络训练样本的数据量很大时，梯度下降法效率很低，应该采用计算效率比较高的随机梯度下降（Stochastic gradient descent）方法或者是迷你批量方法（Mini-Batch）进行优化。深度学习中，神经网络的训练一般采取迷你批量的方式进行优化。迷你批量优化的方式即每次根据部分样本（一个迷你批次内的样本，相对于全体样本集而言只是少量样本）进行梯度方向的计算与迭代优化，所有的迷你批次遍历一次时已经进行了很多次的迭代，因此计算效率比传统的批量梯度下降法要高很多。

公式 (5) 的关键在于偏导数的获取，目前最流行的是 1986 年由 Rumelhart 和 McClelland 提出的反向传播算法（BP 算法），从输出层开始后推，使用误差反向传播的方式对参数进行优化。具体的方法本报告不再赘述。

（五）自编码器和深度网络

深度学习模型事实上是一个含有多个隐层（隐层数量大于等于 2 个）的神经网络模型。原始数据经过一层一层的抽象之后，最后进行分类。但是普通的多层神经网络模型是高度非凸的，在训练时，由于存在大量的局部最优点而且收敛性差，很难获得好的学

习结果，在实际应用时过拟合现象过于严重。而深度学习模型的训练中，先利用大量的数据对无监督网络进行逐层学习，再将训练好之后的参数作为有监督神经网络的参数学习的初始值。模型的总体结构如图 10 所示。该深度学习网络含有两个隐层 H1 和 H2。在逐层的无监督学习中，先通过原始数据 X 学习获得第一个隐含层 H1，然后通过第一个隐含层 H1 学习获得第二个隐含层 H2。当无监督网络训练好之后，一般采用带有标签的样本通过反向传播算法进行有监督学习。

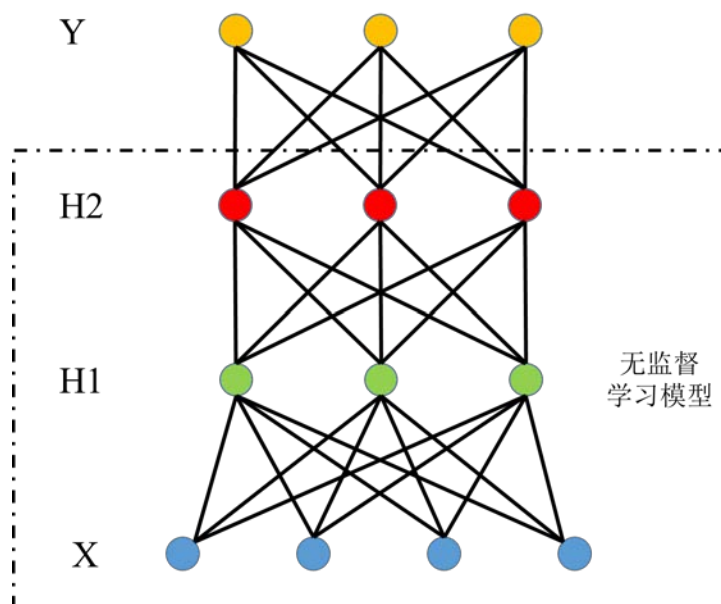
自编码器是一种常用的无监督网络学习方法，也是一种特殊的神经网络模型。该网络模型的输出就是它的输入，如图 11 所示。通过自编码器 $\mathbf{w}^{(1)}$ “编码”获得隐层 Z 之后，可以通过 $\mathbf{w}^{(2)}$ 对隐层进行“解码”，重新还原出原始数据。模型的优化目标函数为

$$L_{AE} = \frac{1}{N} \sum_{i=1}^N \left\| \tilde{\mathbf{x}}^{(i)} - \mathbf{x}^{(i)} \right\|_2^2 \quad (6)$$

即要使得模型的预测输出尽可能等于输入。从信息的传递上来说，这是使得信息损失尽可能少的一种编码方式。

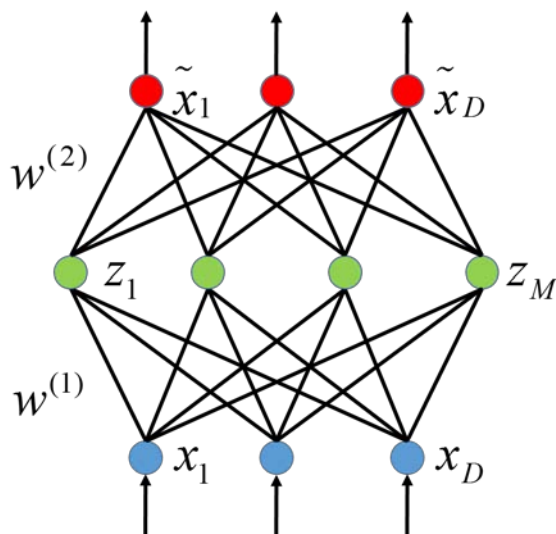
基于自编码器的深度学习模型思路是：对原始的数据，通过自编码器的学习，获得编码层 H1；然后对编码 H1 进行自编码器学习，获得编码层 H2；如此逐层获得深度神经网络的各个隐含层。然后将学习好的编码器系数和结果作为深层神经网络模型的初始值，进行深层神经网络的学习。

图10：深度学习示意图



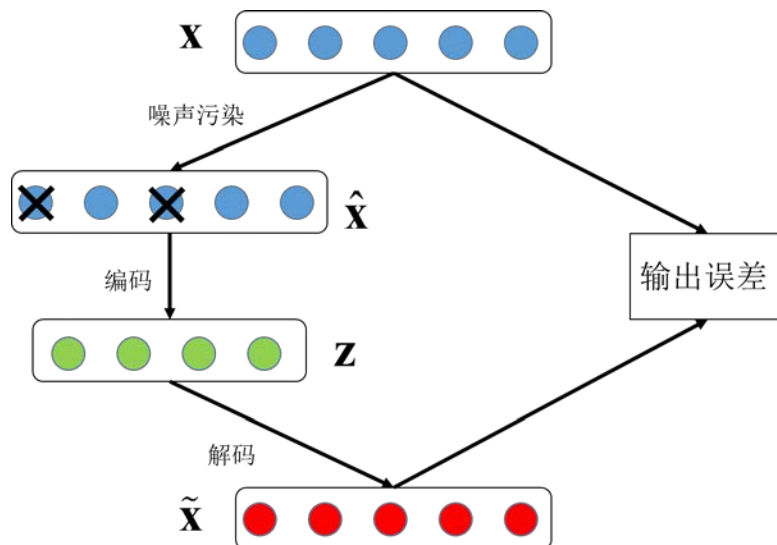
数据来源：广发证券发展研究中心

图11：自编码器示意图



数据来源：广发证券发展研究中心

图12：降噪自编码器示意图



数据来源：广发证券发展研究中心

为了使得自编码器获得的结果更加鲁棒，减少过拟合，Bengio（2008）提出了一种称为降噪自编码器的方式，在每次对自编码器进行训练时，将自编码器输入层的节点以一定概率随机设置为0。如图12所示，其中 $\hat{\mathbf{x}}$ 为将自编码器输入 \mathbf{x} 中节点随机设置为0的结果，可以视为是一种经噪声污染后的形式。在每一次迭代更新参数的时候都需要重新获取 $\hat{\mathbf{x}}$ ，且要使得 $\hat{\mathbf{x}}$ 经过编码解码之后能够重构出尽可能接近原始输入 \mathbf{x} 的结果。实

实践证明通过这种鲁棒自编码器训练获得的深度学习模型比普通的自编码器模型或者受限玻尔兹曼机获得的模型一般要性能更优。本报告采用的就是鲁棒自编码器。

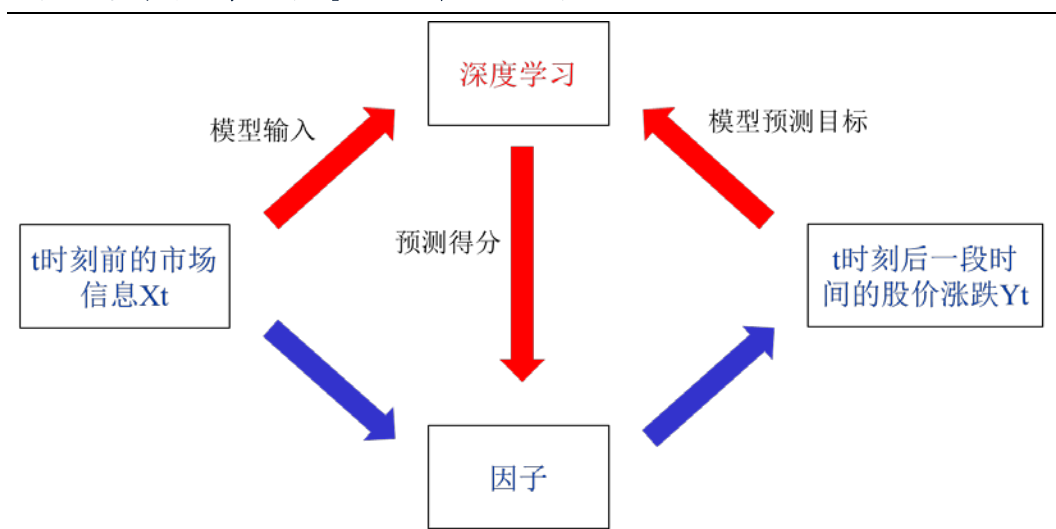
因此，本报告中深度学习模型的训练分成两部分。第一步是使用没有标签或者去掉标签的数据，通过降噪自编码器逐层学习深层网络的隐层系数。第二步是将降噪自编码器学习获得的系数作为网络系数初始值，使用含有标签的数据，通过反向传播算法对深层网络进行有监督学习。

三、实证分析

多因子策略在选取因子时，要使得所选因子在一定的操作周期内具有较大的Alpha，即超额收益。如果某因子与股票未来一段时间内的投资收益直接相关，则该因子可以作为Alpha策略的备选因子。基于深度学习的交易策略就是借助于深度学习对大量的历史交易数据进行学习，建立预测模型，从而获取Alpha因子，即深度学习模型的预测得分。

基于深度学习的Alpha策略如图13所示，深度学习模型建立起当前时刻（t时刻）及此前时刻市场数据 X_t 以及一段时间后股票价格的变化 Y_t 之间的关系，即使用t时刻信息 X_t 通过深度学习模型对此后的 Y_t 进行预测，其预测得分可以作为Alpha策略的因子。

图13：基于深度学习的Alpha因子策略示意图



数据来源：广发证券发展研究中心

（一）深度学习预测模型

首先是深度学习预测模型。一般的，预测时间间隔越短的话，机器学习模型的预测能力会越强。但是短时间内股票价格的涨跌幅度一般较小，考虑到股票交易的成本以及买卖股票对市场造成的冲击，短时期的预测模型很难实现盈利。综合考虑到模型预测的准确程度和交易成本，我们需要在预测周期的选择上进行折中。本报告考虑的是以周为换仓周期的多因子策略，预测模型也选择以周为预测周期，即在每天收盘的时刻进行预测，每次预测的对象都是该交易日之后第5个交易日收盘时刻股价相对于当前收盘价的涨跌情况。

本报告中，多因子策略的选择股票池为中证800成份股（以2014年上半年的标的股票为准），考察时期为自2006年1月以来的股票行情。我们采取滚动预测的方式来建立深度学习模型，即每年训练一次模型，用来预测此后的股票涨跌幅，如图14所示。第一个模型用该股票池内股票在2006年至2010年行情为样本内数据，进行模型训练；在2011年以来的行情为样本外数据。第二个模型用该股票池内股票在2006年至2011年行情为样本内数据，进行模型训练；在2012年以来的行情为样本外数据。第三个模型用该股票池内股票在2006年至2012年行情为样本内数据，进行模型训练；在2013年以来的行情为样本外数据。

图14：基于深度学习的滚动预测模型示意图

1/1/2006	1/1/2011	1/1/2012	1/1/2013
模型2010：样本内		样本外	
模型2011：样本内			样本外
模型2012：样本内			
			样本外

数据来源：广发证券发展研究中心

在Alpha策略中，我们寄望于找到与股票大幅度变化相关的因子。仅有少数周，股票价格的波动会很大。以第三个预测模型为例，样本内102万多个样本中，有22万多个样本的5个交易日后股价变化绝对值大于等于8%，所占比例约20%。因此，预测模型在训练时，选取了三类样本：5个交易日后涨幅大于8%的样本（大涨样本）；5个交易日后跌幅大于8%的样本（大跌样本）；5个交易日后涨跌幅小于1%的样本（平盘样本）。通过深度学习模型，可以对样本是否大涨和大跌给出预测得分。

在预测模型输入的选择上，选择的是短期内的股票价格，以及价格的变化范围，买卖盘价格和委卖委买量等。选取的输入如表2所示。为了在输入数据中引入市场中量价指标随时间的动态变化，我们将该交易日所考察股票的5分钟高频行情序列共48个时间段的行情（ $\mathbf{x}_{t,1}, \mathbf{x}_{t,2}, \mathbf{x}_{t,3}, \dots, \mathbf{x}_{t,48}$ ）加入输入变量序列中，组成扩展的输入向量 $\bar{\mathbf{x}}_t = [\mathbf{x}_{t,1}, \mathbf{x}_{t,2}, \mathbf{x}_{t,3}, \dots, \mathbf{x}_{t,48}]$ （前面50个交易日的收盘价不用再重复选取），因此，每个样本都有386个特征变量（ $7 \times 48 + 50 = 386$ ）。

由于A股市场有涨跌停限制，为了能够使得Alpha策略能够应用于实际交易中并且有足够大的市场容量，本报告在训练预测模型和实证分析选取股票组合时都预先剔除了当天收盘价格相对于上一个交易日收盘价格的涨幅超过9%或者跌幅超过9%的股票。

为了消除不同股票的差异性以及股票处于不同价格时的差异性，数据的预处理过程包括股价数据的标准化（将交易日t的输入向量 $\bar{\mathbf{x}}_t$ 中所有股价数据除以交易日t的收盘价格并取对数），极端数据的平滑（用变量的99.9%和0.1%为门限，超过门限的数据用门限值来替代），和归一化（所有数据都按照公式(7)转化为[0, 1]之间的数据）：

$$x_t^i := \frac{x_t^i - \min x^i}{\max x^i - \min x^i} \quad (7)$$

表2：深度学习股价预测特征选取

选取输入变量	说明
收盘价	
最高价	
最低价	
开盘价	
买卖盘报价平均价格	（买盘报价+卖盘报价）/2
成交量	
委买委卖量之比	$\log(\text{委买量}/\text{委卖量})$
此前50个交易日收盘价格	

数据来源：广发证券发展研究中心

（二）交易策略

本报告中的预测模型和交易策略模型都是基于周频的策略，具体来说，选取每个星期最后一个交易日作为我们建仓平仓的交易日。根据前面的预测模型，我们获得了

个股在T=5个交易日之后大幅上涨预测的得分ScoreUp和大幅下跌预测的得分ScoreDown，由此，我们可以很方便的建立起类似多因子选股模型的交易策略。

具体的来说，以组合规模为100为例，选取当日上涨预测得分ScoreUp最大的前100只股票作为看多组合，当日下跌预测得分ScoreDown最大的前100只股票作为看空组合。由此至少可以衍生出以下三种策略：

(1) 单向做多策略。每个交易日买入前100只股票的看多组合。组合中间的个股按照资金等权的方式进行配置。

(2) 多空配对策略。每个交易日买入前100只股票的看多组合，通过融券卖出前10%的看空组合。看多组合和看空组合总共的资金量相等，组合内部按照资金等权的方式进行个股配置。

(3) 资金等权对冲策略。每个交易日买入前100只股票的看多组合，同时做空沪深300股指期货进行对冲。单次交易盈亏 $=0.8 \times (\text{看多组合收益} + \text{做空股指期货收益})$ 。看多组合内按照资金等权的方式进行个股配置。

以上策略中，单向做多策略风险一般比较大，多空配对策略由于需要融券做空，实际操作上可能会有一些问题，因此本报告中采取第三种交易策略，即资金等权的对冲策略，其中股票组合的交易手续费为单边3‰。

(三) 预测模型效果

深度学习网络的第一个隐层有400个节点，第二个隐层有200个节点，输出层有2个节点（输出ScoreUp和ScoreDown依次表示预测价格是上涨或者下跌的得分）。无监督学习的隐层训练迭代次数为50次，有监督学习的迭代次数为400次。在Intel Xeon E5620，主频2.4GHZ的处理器下，单个预测模型的训练时间约为20小时。模型应用时，单个样本的预测耗时在5ms以内。

样本外预测结果如表3所示，给出了上涨预测得分ScoreUp和下跌预测得分ScoreDown不同时股票价格在5个交易日后的平均变化情况。可以看到，上涨预测得分ScoreUp前5%的样本（平均每个交易日选择购买40只股票）5个交易日后的实际股票价格的平均涨幅为1.97%；ScoreUp前10%的股票（平均每天80只股票）5个交易日后的平均股价涨幅为1.39%，ScoreUp前15%的股票（平均每天120只股票）5个交易日后的平均股价涨幅为1.11%，ScoreUp前20%的股票（平均每天160只股票）5个交易日后的平均股价涨幅为0.92%。

表3: 深度学习股价预测结果

预测得分	该样本 5 个交易日后股价变化均值			
	前 5%	前 10%	前 15%	前 20%
ScoreUp	1.97%	1.39%	1.11%	0.92%
ScoreDown	-0.66%	-0.45%	-0.33%	-0.26%
ScoreUp-ScoreDown	1.97%	1.48%	1.24%	1.09%
ScoreDown-ScoreUp	-0.64%	-0.59%	-0.56%	-0.50%

数据来源: 广发证券发展研究中心, 天软科技

(四) 对冲策略效果

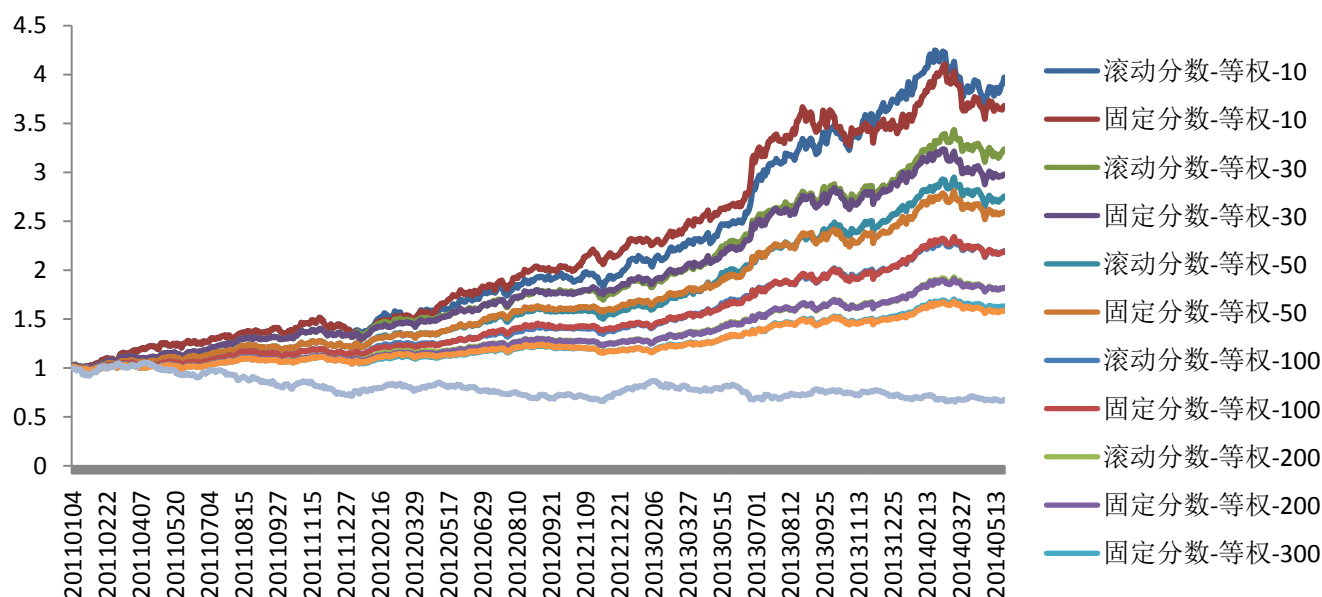
为了确保模型的真实可靠性, 本报告采取了模拟实时追踪的方式对模型进行检验。也即确定好建仓日之后分析某一建仓日的盈亏时, 只从天软数据库读入截止到该交易日收盘时刻为止的数据, 进行深度学习模型预测。根据预测模型的得分选择要买入的个股组合, 再回到天软读取一周后的数据进行该次交易的盈亏计算。

我们依次考虑了不同的组合规模, 即组合内股票数量为 10 只, 30 只, 50 只, 100 只, 200 只, 300 只时对冲策略的实证效果。为了说明滚动预测模型的效果, 比较了滚动预测模型和固定模型 (仅用 2010 年之前的数据建立一个深度学习预测模型, 所有样本外数据都用这个预测模型进行预测) 下的收益情况。其收益曲线如图 15 所示, 按年度的收益率如表 4 所示。

可以看到, 当组合规模越小时, 累积收益率越高, 但是回撤也会相应增大。这是因为, 当上涨预测得分 ScoreUp 越大时, 一般会对应更大的上涨可能性。对冲后回撤依然比较大是因为当组合内股票数量较少时, 沪深 300 股指期货对冲的效果不明显。一般而言, 滚动预测模型的效果会稍微优于固定好的预测模型。

随着组合规模的放大, 因子效应缩减, 累积收益率降低, 但是回撤也相应减少。综合考虑收益率及其稳定性, 组合规模控制在 100 左右比较合适。

图15: 不同组合规模的股票等权组合策略收益曲线



数据来源: 广发证券发展研究中心, 天软科技

一月成交金额和股价动量是主流多因子模型中性能比较好的两个量价因子, 基于成交金额-股价动量双因子的对冲组合与深度学习预测模型对冲组合的比较如图 16 所示。两个组合都是组合规模为 100 的等权组合, 按周换仓。可以看到, 深度学习对冲组合的收益情况大大优于成交金额-股价动量双因子对冲组合。按月的收益比较如图 17 所示, 在绝大多数月份, 深度学习对冲组合的收益都好于成交金额-股价动量双因子组合。

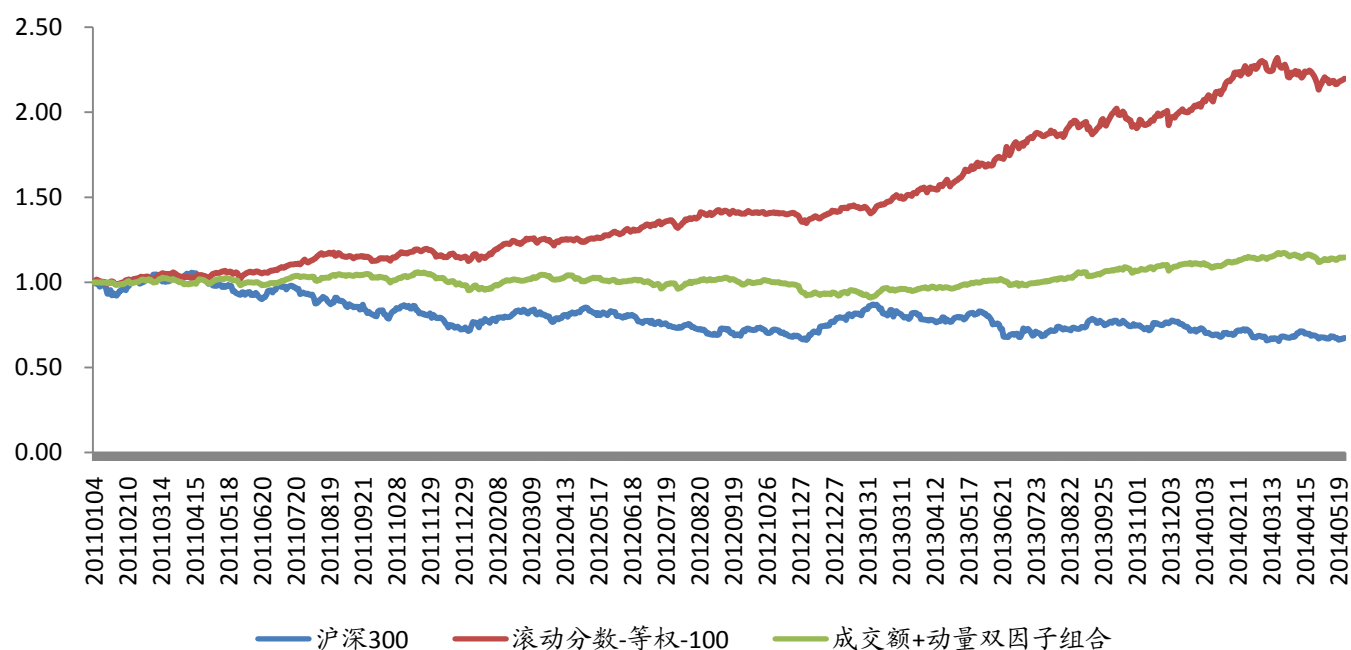
表4: 不同数量股票等权组合策略年度收益率

	全样本	2011	2012	2013	2014
滚动分数-等权-10	297.20%	39.30%	47.10%	82.00%	5.00%
固定分数-等权-10	268.10%	39.30%	62.30%	50.30%	7.40%
滚动分数-等权-30	223.70%	34.30%	36.50%	58.40%	9.60%
固定分数-等权-30	197.90%	34.30%	38.30%	53.90%	2.90%
滚动分数-等权-50	175.60%	22.80%	30.10%	61.30%	5.80%
固定分数-等权-50	159.60%	22.80%	34.40%	48.40%	5.20%
滚动分数-等权-100	119.80%	15.20%	23.40%	43.50%	6.80%
固定分数-等权-100	119.10%	15.20%	25.30%	41.70%	6.60%

滚动分数-等权-200	82.10%	9.10%	17.10%	32.80%	7.00%
固定分数-等权-200	81.90%	9.10%	16.90%	33.60%	6.60%
固定分数-等权-300	63.10%	7.20%	11.10%	29.00%	6.20%
滚动分数-等权-300	58.70%	7.20%	10.90%	27.40%	4.60%

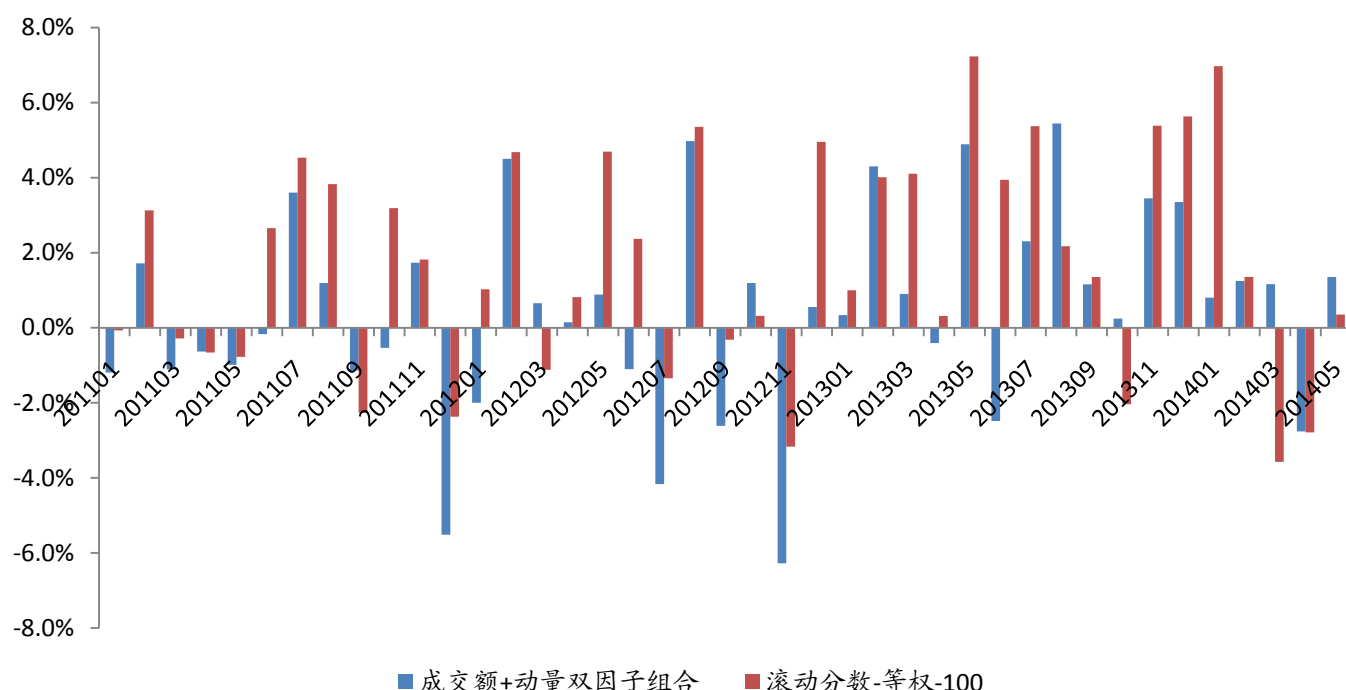
数据来源：广发证券发展研究中心，天软科技

图16：深度学习对冲组合策略与成交金额-股价动量双因子对冲组合收益曲线对比



数据来源：广发证券发展研究中心，天软科技

图17：深度学习对冲组合策略与成交金额-股价动量双因子对冲组合月度收益率对比

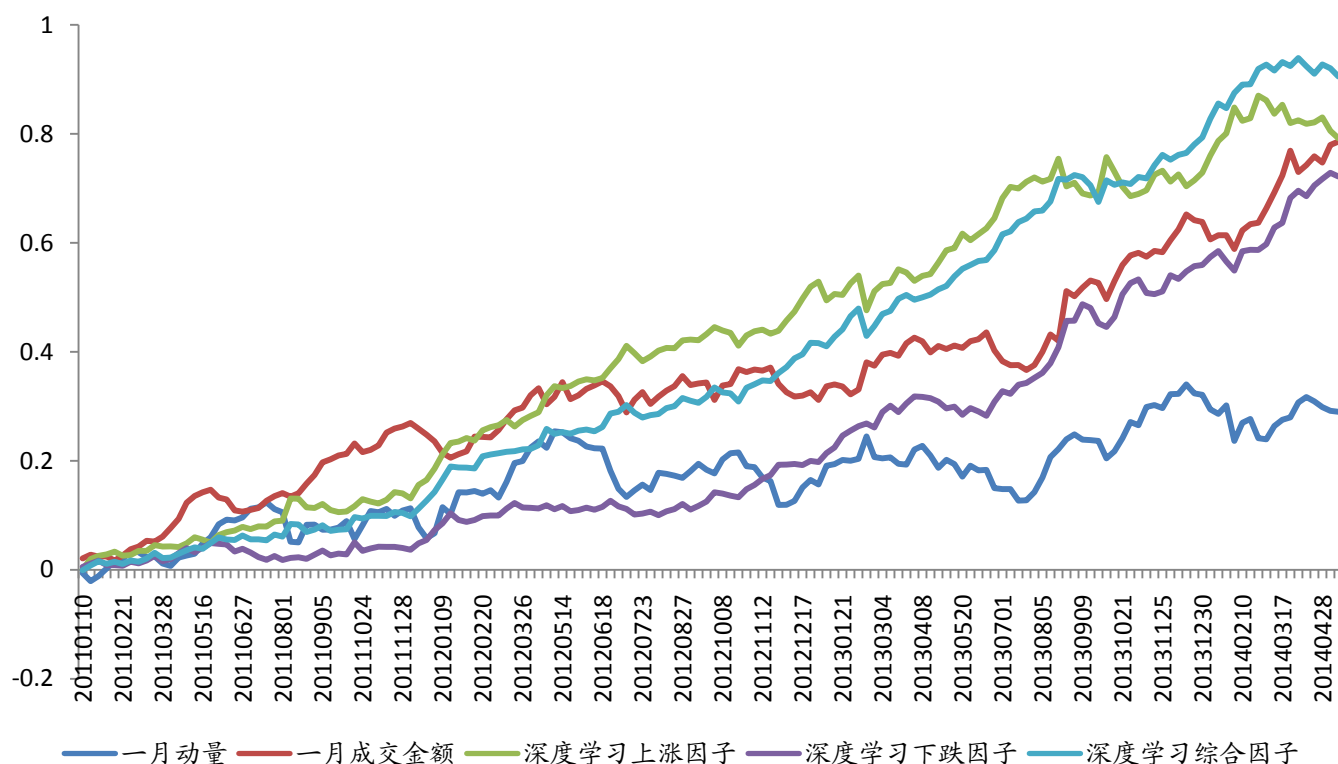


数据来源：广发证券发展研究中心，天软科技

除了将上涨预测得分 $ScoreUp$ 作为选股 Alpha 因子外，也可以将下跌预测得分 $ScoreDown$ 作为选股因子，或者将两个预测得分之差 $ScoreUp - ScoreDown$ 作为选股的综合因子。考虑这些因子各自的五档多空累计表现，即每一期组合的股票选取时做多该因子比较大的前 160 只股票，做空后 160 只股票。各自的累积收益曲线如图 18 所示。各因子的 IC, IR 等表现如表 5 所示。可以看到，深度学习策略的股票组合换手率比较高，以周为持仓周期的多因子策略，换手率都在 40% 以上。从因子信息比率 IR 上来说，深度学习预测模型各因子都具有比主流多因子模型更大的 IR。

从行业分布来看，图 19 给出了滚动预测模型组合规模为 100 时每期组合所选行业股票数量的变化情况。从中也可以看到该策略组合的股票换手率比较高，每一期所选择的行业与其上一期的组合也有很大的变化。图 20 展示了每期组合各行业股票数量中位数。医药，房地产，电器元器件等行业的股票是该策略比较青睐的股票。

图18：深度学习不同预测因子及其组合的五档多空累计表现对比



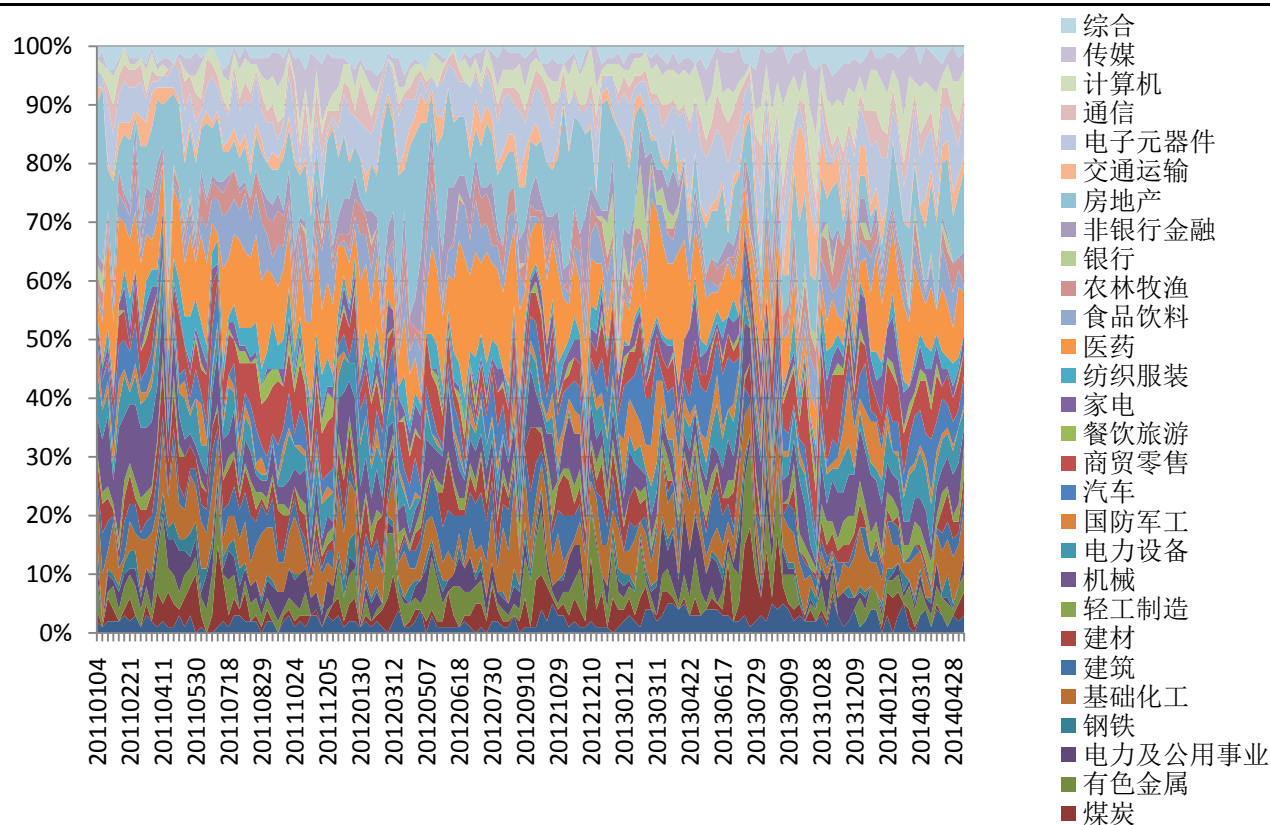
数据来源：广发证券发展研究中心，天软科技

表5：深度学习不同因子的表现对比

因子名称	IC	IR	胜率	换股比例
一月动量	-1.90%	0.75	54%	38%
一月成交金额	-3.51%	2.03	66%	11%
深度学习上涨因子	3.92%	2.34	67%	42%
深度学习下跌因子	-3.92%	2.62	62%	52%
深度学习综合因子	4.60%	3.23	68%	52%

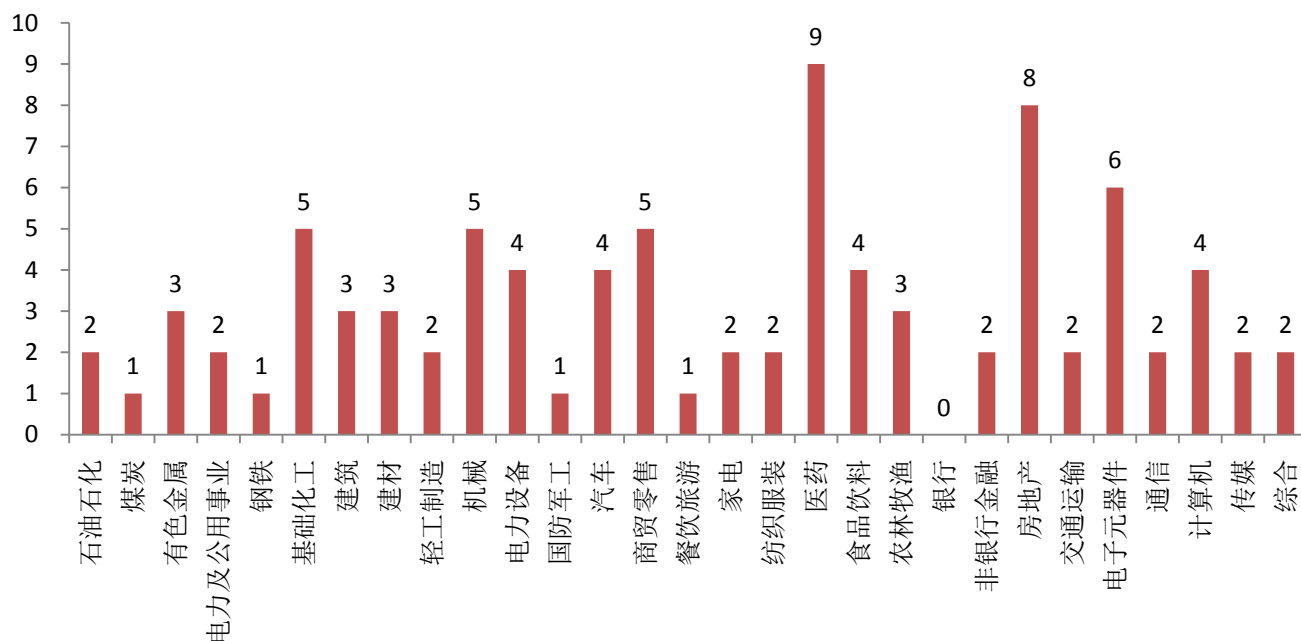
数据来源：广发证券发展研究中心，天软科技

图19：深度学习Alpha策略股票组合行业分布图



数据来源：广发证券发展研究中心，天软科技

图20: 深度学习Alpha策略样本内各期股票组合行业数量中位数



数据来源：广发证券发展研究中心，天软科技

四、总结与讨论

本报告通过深度学习算法对股票市场数据进行挖掘,建立起通过股票市场数据预测股价短期内走势的模型,通过该预测模型的预测得分,我们可以筛选出股票组合并且获得超额收益。因此,该预测得分是可以产生 Alpha 收益的有效因子。与传统的量价因子相比较,通过深度学习算法的因子能够获得更好的收益率。

风险提示

策略模型并非百分百有效,市场结构及交易行为的改变以及类似交易参与者的增多有可能使得策略失效。

广发证券—行业投资评级说明

- 买入： 预期未来 12 个月内，股价表现强于大盘 10%以上。
- 持有： 预期未来 12 个月内，股价相对大盘的变动幅度介于-10%~+10%。
- 卖出： 预期未来 12 个月内，股价表现弱于大盘 10%以上。

广发证券—公司投资评级说明

- 买入： 预期未来 12 个月内，股价表现强于大盘 15%以上。
- 谨慎增持： 预期未来 12 个月内，股价表现强于大盘 5%-15%。
- 持有： 预期未来 12 个月内，股价相对大盘的变动幅度介于-5%~+5%。
- 卖出： 预期未来 12 个月内，股价表现弱于大盘 5%以上。

联系我们

	广州市	深圳市	北京市	上海市
地址	广州市天河北路 183 号 大都会广场 5 楼	深圳市福田区金田路 4018 号安联大厦 15 楼 A 座 03-04	北京市西城区月坛北街 2 号 月坛大厦 18 层	上海市浦东新区富城路 99 号 震旦大厦 18 楼
邮政编码	510075	518026	100045	200120
客服邮箱	gfyf@gf.com.cn			
服务热线	020-87555888-8612			

免责声明

广发证券股份有限公司具备证券投资咨询业务资格。本报告只发送给广发证券重点客户，不对外公开发布。

本报告所载资料的来源及观点的出处皆被广发证券股份有限公司认为可靠，但广发证券不对其准确性或完整性做出任何保证。报告内容仅供参考，报告中的信息或所表达观点不构成所涉证券买卖的出价或询价。广发证券不对因使用本报告的内容而引致的损失承担任何责任，除非法律法规有明确规定。客户不应以本报告取代其独立判断或仅根据本报告做出决策。

广发证券可发出其它与本报告所载信息不一致及有不同结论的报告。本报告反映研究人员的不同观点、见解及分析方法，并不代表广发证券或其附属机构的立场。报告所载资料、意见及推测仅反映研究人员于发出本报告当日的判断，可随时更改且不予通告。

本报告旨在发送给广发证券的特定客户及其它专业人士。未经广发证券事先书面许可，任何机构或个人不得以任何形式翻版、复制、刊登、转载和引用，否则由此造成的一切不良后果及法律责任由私自翻版、复制、刊登、转载和引用者承担。