上野野野野生物化学与生物物理进展 Progress in Biochemistry and Biophysics 2016, 43(5): 472~483 www.pibb.ac.cn

生物医学数据分析中的深度学习方法应用*

李 渊 ¹⁾ 骆志刚 ^{1)**} 管乃洋 ¹⁾ 尹晓尧 ¹⁾ 王 兵 ³⁾ 伯晓晨 ²⁾ 李 非 ^{2)**} (¹⁾国防科学技术大学计算机学院软件研究所,长沙 410073; ³⁾军事医学科学院放射与辐射医学研究所,北京 100850; ³⁾中国人民解放军 63928 部队,北京 100028)

摘要 生物医学数据的积累速度史无前例,为生物医学研究带来机遇的同时,也让传统数据分析技术面临巨大挑战.本文综述了深度学习方法应用在生物医学数据分析中的最新研究进展.首先阐述了深度学习方法,列举深度学习方法的主要实现模型,随后总结了目前生物医学数据分析中的深度学习方法应用情况,分析了在数据处理、模型构建和训练方法等方面共有问题的解决方法,最后给出了深度学习方法应用于生物医学数据分析时可能存在的问题及建议.

关键词 深度学习,高通量组学,临床医学,数据挖掘学科分类号 Q811.4

DOI: 10.16476/j.pibb.2015.0339

自 1970 年 Sanger 测序法出现,随着高通量组学技术(high-throughput omics technologies)的蓬勃发展,生物医学正加速进入大数据时代. 其中以基因组为代表的组学数据的积累速度史无前例. 近年著名的生物数据项目有千人基因组计划(1000Genomes) [1]、DNA元件百科全书计划(ENCODE) [2]、表观组学路线图计划(Roadmap Epigenomics Program) [3]、细胞印记整合网络数据(LINCS) [4]、基因表达数据库(GEO) [5]、癌症基因组图谱计划(TCGA) [6]等.

借助生物医学多组学、多层次的海量数据,人们可以从前所未有的广度和深度来研究生物体运行机制,这对深入探索人类疾病的分子机理、推动个体化精准医疗[^{7-8]}有重大意义.但海量复杂数据带来的挑战同样巨大:一方面要处理的生物医学数据维数更高、要求模型有更强的学习适应能力;另一方面生物医学大数据更加分散破碎,数据结构复杂,常常要整合不同类型的信息^[9],如基因组、蛋白质组和临床研究的数据等.这些问题都对数据特征提取提出了更高的要求.因此在生物医学大数据时代,统计分析方法及数据挖掘方法都面临巨大的挑战.

近年来,深度学习方法[10]得到广泛应用,已经在图像识别[11-12]、语音识别[13-14]等领域取得了令人

瞩目的成果. 作为机器学习中的重要方法之一,由 于其强大的自动特征提取、复杂模型构建以及图像 处理能力,非常适合处理生物医学数据分析所面临 的新问题, 引起了生物医学领域研究人员的广泛关 注. 深度学习方法从人工神经网络模型发展而来. 通过组合多个非线性处理层对原始数据进行逐层抽 象,从数据中获得不同层面的抽象特征并用于分类 预测. 与传统机器学习方法相比, 具有以下三个特 点: a. "深层"模型架构. 深度学习模型的多层结 构与动物的视觉处理系统极为相似[15]. 与其他浅层 模型,如支持向量机(support vector machine, SVM)等相比,深度学习模型拥有更多的隐层,包 含更多的非线性变换,这使得深度学习拟合复杂模 型的能力大大增强. b. 多层数据特征表示[16]. 深 度学习模型以数据的原始形式作为输入,之后将当 前层的输出作为下一层的输入,逐层堆叠,由此归

^{*}国家高技术研究发展计划(2015AA020100), 国家自然科学基金重大 计划 (U1435222), 国家自然科学基金面上项目(81273488, 61402486)和传染病防控关键技术研究项目(BWS14C051)资助.

^{**} 通讯联系人.

骆志刚. Tel: 0731-84575835, E-mail: zgluo@nudt.edu.cn 李 非. Tel: 010-66932251, E-mail: pittacus@gmail.com 收稿日期: 2016-01-21, 接受日期: 2016-03-25

纳得到更高级的特征表示,从而能够刻画复杂数据结构.c.无监督学习.深度学习模型在训练中加入无监督学习过程,通过预训练获得良好的模型初值,能有效提升训练效果,另外无标签数据加入训练也增加了可用数据的规模.

本文首先介绍深度学习中最重要的几种模型及深度学习模型的构建和训练过程,再分别列举介绍医疗数据和生物数据分析中现有的深度学习应用,之后总结分析现有应用在数据处理、模型构建和训练方法上的共同点及遇到的问题,最后分析未来生物医学领域深度学习可能的应用场景.

1 深度学习方法

深度学习方法包含多种深度模型,其中通用模型以深度信念网络模型(deep belief networks, DBN) [17] 和堆叠自动编码器模型(stacked auto-encoder,SAE)[18]为代表,另外有用于图像处理的卷积神经网络模型(convolution neural nets,CNN)[19]和用于序列数据处理的循环神经网络模型(recurrent neural nets,RNN)[20-23]. 除此之外,近年来还出现了多种新的深度模型,如由经典模型衍生而来的随机生成网络模型(generative stochastic Network,GSN)[24]、基于独立子空间分析网络(independent subspace analysis network,ISA)形成的堆叠网络模型(stacked SPN)[26]等.

上述模型中, DBN、SAE 和 CNN 模型在生物 医学数据分析中应用最为广泛,并且其模型结构、训练过程等具有一定的代表性. 本节以这三种深度 学习模型为例,详细介绍深度模型的基本单元、构建过程和训练过程,如图 1 所示.

1.1 深度信念网络 DBN

DBN于 2006 年由 Hinton 等¹¹⁷提出,由多个限制波尔兹曼机(restricted boltzmann machine, RBM)作为基本单元堆叠而成.单个 RBM 包含可视层和隐藏层,两层之间双向连接,其中可视层同时作为输入输出复用,其拓扑图见图 1a. RBM 模型使用对比散度算法(contrastive divergence)对无标记样本进行训练,属于无监督学习算法.最终训练完成的模型能够从隐藏层数据反向还原可视层数据,即隐藏层是可视层数据的抽象表达.

DBN 构建过程如下: 首先训练得到第一个RBM, 随后冻结模型的权值并将其隐藏层作为下

一个 RBM 模型的可视层,用同样的方法可训练得到第二个 RBM. 依次类推,可以得到多个 RBM. 将多个 RBM 按顺序堆叠在一起便构成一个深度玻尔兹曼机(deep boltzmann machine,DBM)模型. 此时模型的输出将是输入数据经过多次抽象后得到的多层抽象表示,也就是模型自动学习到的数据特征. 若将此特征作用于分类器,通常能得到好的分类结果. 在 DBM 的顶端加入"联想记忆"层,则构成 DBN 模型. 如果在 DBN 第一层之前加入卷积处理层,可得到卷积深度信念网络模型(convolution DBN). 该模型已成功应用在人脸识别[27]、音频分类[28]问题中.

1.2 堆叠自动编码器 SAE

SAE 由 Bengio 等[18]提出,其基本元件是自动编码器(Auto-encoder,AE). AE 包含输入层、隐藏层和输出层,三层之间逐级连接,其拓扑图见图 1a. AE 模型将训练目标设为拟合输入数据,即设定网络输出等于输入,随后使用反向传播算法训练. 虽然 AE 模型的训练过程基于有监督学习算法,但并不要求原始数据有分类标签,因此整个训练过程仍是一个无监督学习过程. 若在训练中加入稀疏惩罚项,即对网络中被激活单元的个数加以限制,便构成稀疏自动编码器(sparse AE) ,如果训练中在输入数据中加入随机噪声,便构成去噪自动编码器(denoising AE)[29],这两种模型在实际中往往能学到更好的数据特征.

SAE 的构建过程与 DBN 类似,模型构建过程 见图 1b. 训练得到第一个 AE 后,将其隐藏层作为输入,用同样的方法可训练第二个 AE,依次类 推可训练得到多个 AE. 依次将多个 AE 堆叠在一起,便构成 SAE 模型,此时 SAE 的最后一层是输入数据经过多次变换处理后得到的抽象特征. 最后再根据问题不同设定,连接不同的输出层,通过有监督学习算法训练输出层的权值,从而得到最终分类结果.

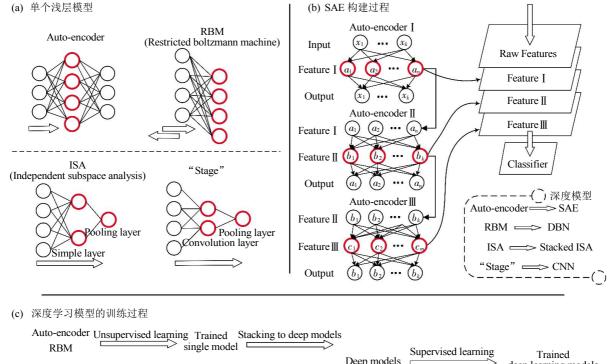
1.3 卷积神经网络 CNN

CNN 由 Lecun 等¹¹⁹提出,主要用于图像处理、图像识别等领域,如大规模图像识别深度学习网络GoogLeNet^[30]和 Adam^[31]等. CNN 的基本思想源于对猫视觉神经^[32]的研究,其中感受野(receptive field)原理的发现对 CNN 模型意义重大. CNN 的关键思想在于多层堆叠、区域连接、权值共享和池化(pooling).

CNN 模型由多个"stage"堆叠而成,每一个

基本元件"stage"的拓扑结构如图 la 所示,它包 含一个卷积层和一个 polling 层, 卷积层能捕捉图 片中的区域性连接特征,且应用了权值共享原理, 使模型要训练的参数个数大大减少. pooling 层将 相邻的多个节点合并为一个来合并相似特征,进一

步减小训练的数据量. 构建多个"stage"后,将它 们堆叠在一起, 在模型末端加入多个全连接层和分 类器便构成了 CNN, 最后用有监督方式对 CNN 整 体进行训练.



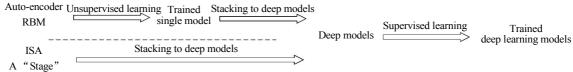


Fig. 1 Architecture, construction process & training process of deep learning

图 1 深度学习模型的基本架构、构建过程和训练过程

(a) 4 个子图分别为组成 SAE、DBN、Stacked ISA 和 CNN 4 个深度学习模型的基本单元. (b) 由 AE 构建 SAE 模型的过程,其他深度学习模 型的构建过程与之相似. (c) SAE、DBN、Stacked ISA 和 CNN 4 个深度学习模型的训练过程,可基本代表现有深度学习模型的训练过程.

2 生物医学数据分析中的深度学习应用

在可预期的将来,深度学习方法将在生物医学 研究领域中得到越来越广泛的应用. 近年已有许多 研究团队尝试将深度学习方法应用在生物医学数据 分析处理中,为进一步的研究工作提供了重要的指 引. 本节分别从医疗数据和生物数据两个方面,简 要介绍近年来深度学习应用方面的研究进展,具体 工作见表 1.

2.1 医疗数据分析中的深度学习应用

2.1.1 疾病诊断

疾病诊断是深度学习在医学上的主要应用之

一. 它基于患者的疾病相关数据,通过深度学习模 型预测异常病变或发病风险,进行疾病的辅助诊 断. 自动化的疾病辅助诊断能更快地处理数据,能 为医师提供参考, 且其判断不易受到主观因素的干 扰,在减轻医师工作负担的同时提升效率和诊断准 确率. 自动疾病诊断包括疾病诊断、疾病分类和疾 病分级等方面.

2011年, 宾夕法尼亚大学 Wulsin 等[3]使用 DBN 对脑电波波形图建模,进行人类脑部异常检 测,此方法比 SVM 拥有更快的速度,从而具有更 好的实时性. 2014 年, Chakdar 等[34]使用 DBN 进 行基于子宫抹片检查的低级别鳞状上皮内病变(low grade squamous intraepithelial lesion, LGSIL)诊断, 该方法能从抹片图像中自动提取特征进行疾病诊断,实验中, DBN 提取出的特征和原始特征共同用于 SVM 模型,可使分类准确率达到 100%.

2015年,新加坡 Gao 等[55]结合使用 CNN 和递归神经网络(recursive NN),基于眼部检查图像对核性白内障进行严重程度分级,深度学习方法打破了该领域之前的记录. 2014年,约翰霍普金斯大学的 Yang 等[56]采用 SAE 模型对脑部核磁共振图像建模,进行小脑运动失调症的分类,其分类准确率可

达 97%.

此外值得一提的是深度学习在阿尔茨海默病 (Alzheimer disease, AD)研究中的大量应用. 轻度 认知障碍 (mild cognitive impairment, MCI)是 AD 的早期症状,但并非所有的 MCI 都会发展成 AD,因此判断 MRI 患者的类型及预测医疗干预的效果十分重要. 近几年来,美韩两国的多家研究机构都使用深度学习模型,基于核磁共振成像 MRI 和正电子发射断层扫描图 PET,进行 AD/MCI 分类研究[37-41],其分类效果十分突出.

Table 1 Applications of deep learning in biological and medical data analysis 表 1 生物医学数据分析中的深度学习应用

	衣 1 生物医子数据方机中的床皮子习应用					
类别	时间	问题	所用模型			
疾病诊断	2014~2015	AD/MCI 分类[38, 40-44]	SAE/DBN			
	2013~2014	癌症 / 肿瘤诊断[34, 45-47]	CNN/SAE/DBN+SVM			
	2011、2014	脑部疾病诊断[33,36,48]	SAE/CNN/DBN			
	2015	核性白内障分级[49]	CNN+SVM			
	2014	慢性胃炎诊断[50]	DBN			
医学图像处理	2013~2014	医学图像自动分割[51-53]	Stacked ISA/CNN			
	2013	图像关键点发现[46,54-55]	CNN/DBN/CNN			
	2015	MRI 图像重构 ^[56]	SAE			
医学数据建模	2013~2015	脑发育、脑回路建模[57-62]	DBN/CNN			
	2013	情绪分析 ^[©]	DBN			
蛋白质结构预测	2014~2015	二级结构预测[64-66]	SAE/DBN/SAE			
	2012、2014	三级结构预测[67-70]	DBN/3D NN/GSN/DST-NN			
	2015	基于模板的结构预测问	DBN			
测序数据处理	2015	DNA 标注 ^[72]	DNN(Deep NN)			
	2014~2015	RNA 可变剪切分析[73-74]	SAE+DNN			
	2015	非编码 RNA 发现[^{75]}	lncRNA-MFDL			
表达谱数据分析	2013~2015	基于表达谱的亚型鉴定[45,76-77]	DBN/SAE			
其他	2015	药物设计与发现™	DNN			
	2014	蛋白质模型质量评估[79]	SAE			

2.1.2 医学图像处理

医疗机构的医学图像产出量巨大,图像数据往往包含大量潜在信息.目前主要依靠人工判读分析,效率较低且能挖掘到的信息有限,无法充分利用数据资源.深度学习在图像处理领域的优异表现为医学图像的自动化处理提供了新方法.目前深度学习在医学图像中主要应用于临床图像分类,关键目标发现和图片自动分割等方面.

2013 年, Cruz-Roa 等[46]将深度模型用于图像中肿瘤细胞的自动发现,该模型的准确率相比传统

方法有 7%的提升,对癌症自动诊断有重要意义;同年,Carneiro 等[54]使用定制深度模型从超声波数据中追踪左心室心内膜,在超声波数据的自动分析应用方面取得良好的结果. 2012 年,瑞士 Cireşan等[55]将 CNN 用于乳腺癌细胞图片中有丝分裂的自动寻找,该模型的准确率远远超过了以往方法,赢得 当 年 ICPR (international conference pattern recognition)竞赛的冠军.

医学图像分割是医学图像处理的基础,对后续 的病症定量分析、组织三维可视化及治疗方案的制 定都十分重要. 2013 年,芝加哥大学 Aytekin Oto 等[51]使用 Stacked ISA 模型进行前列腺核磁共振图像的自动分割,此模型能够自动从核磁共振图中分割出前列腺部分,使用深度学习方法获得的抽象特征代替以往手工设计的特征,明显提升了图片分割准确率. 同年,哥本哈根大学 Prasoon 等[52]结合使用 3 个二维 CNN 来处理膝盖软骨组织的三维图像,并对其自动分片,该模型的实验结果超过了直接使用图像三维特征的模型. 2014 年,Song 等[53]使用 CNN 模型进行宫颈细胞图片的质核分离,达到 91.34%的准确率.

2.1.3 医学数据建模

深度学习也被应用于医学数据建模.相比疾病诊断,建立模型的问题更加困难但更有意义,其处理对象多为复杂结构或复杂过程,好的模型会有更加广泛的应用.例如疾病发生发展过程模型对相关疾病的分析、监测及预防等都有帮助.

2013 年,密苏里大学 Wang 等[63]用 DBN 对生理数据建模来预测人的情绪,该模型以普通人体体征数据为输入,情绪预测准确率与采用专家设计特征的浅层模型相当,但无需处理的数据作为输入使其更适合大数据问题. 2014 年,Nguyen 等[79]基于SAE 设计了 DL-PRO 模型对蛋白质结构预测模型的质量进行评估,其实验效果超越了该领域之前的最优结果.

深度学习还广泛用于脑部问题建模. 2014年,Brosch等[58]又结合使用三个DBN对大脑形态变化建模,该模型能自动捕捉脑部病变前兆,对脑部疾病的预测预防有重要意义. 同年,Zhang等[59]使用深度CNN,以多模态核磁共振(MR)图像作为输入,对婴儿脑部发育图片中的灰质和白质自动切割,该模型还能根据图片信息区分婴儿大脑的发育阶段.

2.2 生物数据分析中的深度学习应用

相比医学问题,生物学中问题模型复杂,数据量更庞大.目前,深度学习技术主要被用在蛋白质结构预测、测序数据处理和表达谱数据处理三个方面.深度学习在生物数据分析中的优秀表现为探究生物序列和分子结构与疾病的关系提供了新技术.

2.2.1 蛋白质结构预测

蛋白质结构预测极其重要却又十分复杂,基因组测序的完成揭示了大量蛋白质的氨基酸序列,但如何从氨基酸序列推导出蛋白质的高级结构却仍然没有取得明显突破.蛋白质的高级结构对深入研究蛋白质的理化性质和作用机理十分重要.目前,深

度学习模型主要被用于间接预测蛋白质高级结构, 典型的方法有通过氨基酸序列预测蛋白质的二级结构、三级结构或预测响应蛋白质所属模板.

对二级结构预测,2015年,澳大利亚 Lyons 等[64]使用 SAE 模型,通过氨基酸序列预测其主干角度和二面角,该模型表现十分优秀。同年,Spencer等[65]综合使用三个 DBN 模型预测蛋白质的二级结构,该领域研究者之前普遍认为 80%是这一问题的准确率极限,但 Spencer 等的模型准确率达到了 80.7%.

对三级结构预测,2012 年,美国 Di Lena 等[67] 使用多个深度模型,通过三个阶段的优化,来解决蛋白质三级结构预测问题,该模型准确率达到30%,之前该领域的普遍结果仅有20%准确率.之后,Di Lena 等[68]又将时间概念引入蛋白质结构预测,通过在空间和时间两个维度上构建深度模型来预测蛋白质三级结构,相比传统方法有30%的提升.2014 年,普林斯顿大学 Zhou 等[69]对基于 GSN模型做出修改,将其由生成模型改为区分模型并加入卷积处理过程,来处理蛋白质三级结构预测问题,该模型在同一个测试集上的表现比之前的最好结果准确率高出2%.

模板法预测蛋白质结构是通过比较新蛋白与一系列已经被研究清楚的"模板"蛋白的相似性,来预测新蛋白质结构的方法,该方法目前被广泛使用. 2015 年,Taeho等四使用 DBN 模型预测两个氨基酸序列是否对应同一"模板",并依此来预测蛋白质结构,准确率最高可达 91.2%.

2.2.2 测序数据处理

2014年,多伦多大学 Frey 等[^{73-74]}通过构建深度神经网络,分析 mRNA 的组织特异性剪切模式及可变剪切与人类疾病的关系,该模型的表现优于同一团队之前采用贝叶斯网络的结果. 2014年,Fan 等^[75]基于深度学习算法构建了 IncRNA-MFDL,用来自动识别非编码 RNA,为高通量测序成果的后续分析提供了更加有力的方法.

2015 年,Quang 等[^{72]}建立了一个五层深度网络,将基因突变的相关特征作为输入,来判断它是否为致病基因,该模型对编码区和非编码区的序列数据均有效,相较之前的 SVM 模型,其 AUC 有0.09 的提升. 2014 年,加拿大滑铁卢大学 Ibrahim等^[77]通过构建深度模型,根据基因表达数据来寻找对疾病预测具有最大区分度的基因,相比传统方法,在不同病症上有6%~10%的准确度提升.

2.2.3 表达谱数据分析

2013 年,Fakoor 等[45]使用 SAE 模型通过分析基因表达数据进行癌症诊断,该模型不仅对不同组织器官的数据均适用,且能够通过训练获得不同癌症数据的通用特征. 2014 年,加拿大滑铁卢大学Ibrahim 等[77]通过构建深度模型,根据基因表达数据来寻找对疾病预测具有最大区分度的基因,相比传统方法,在不同病症上有 6%~10%的准确度提升. 2015 年,Liang 等[76]设计了多模型 DBN 用于不同平台癌症数据的聚类,该模型不仅能对卵巢癌和肺癌实现准确的亚型分类,且发现了 miR-29a 的表达水平与癌症患者存活时间之间的密切关系.

3 深度学习模型构建及训练

3.1 模型构建

选择合适的深度模型很重要.在现有工作中,SAE 和 DBN 模型应用最广泛,因为它们能很容易地被应用在不同类型问题和数据上,且能保证得到较好的结果. CNN 模型主要被用于图片数据,另外也有研究者选择最新的深度模型(如 GSN 等).此外,组合不同类型浅层模型来构建特别的深度模型有时也会有出人意料的效果. 例如有研究者将卷积处理过程加入 DBN 模型[27-28]使其可用于图像处理,而 Wu 等[60]的工作则是在 CNN 之前加入非监督学习过程,使模型能自动学习数据特征. 另外有研究者将以前使用的浅层模型堆叠在一起构成新的深度模型,如 Stacked SPN[26]、Stacked ISA[25]等.

对某些复杂的生物医学问题,单个深度模型无法直接使用或者训练效果不佳.这时可以组合多个深度模型来构建更庞大、深层、复杂的深度模型解决方案. Heffernan 等[66]通过构建多个深度模型,同时预测蛋白质结构相关的三个参数并基于上一轮的结果不断迭代,最终得到较好的预测结果;Prasoon等[52]则对三维图像分别从x-y、x-z、y-z 三个维度构建 CNN 模型并相互结合来进行三维图像分析,其表现优于直接使用图像三维特征建模. 另外,通过引入不同维度的空间和时间数据,将深度学习模型在多个维度上组合可构建更复杂有效的模型,如空间上三维加时间上一维的模型可用于气象预报、机器人运动轨迹预测等,这种模型构建方式在生物医学数据分析中也大有可为.

实际应用中,若问题模式明确且并无明显特征,可考虑先采用经典深度学习模型进行学习,之后根据问题演变及训练效果调整模型.若已经明确

为图像数据分析,则最好使用 CNN 模型. 若研究者在之前工作中有熟悉的浅层模型,则可尝试搭建自己的深度学习模型. 总之,如果研究者面对的问题较复杂,那么应先理清思路,将目标问题拆分为多个适合深度模型处理的小问题,之后通过在多个维度上组合不同的深度学习模型来进行研究.

3.2 模型训练

深度学习由神经网络模型发展而来,其训练方法也继承了原有训练方法,即反向传播训练方法和随机梯度下降方法,其主要不同在于无监督学习作为预训练的广泛使用. 2010 年,Martens 等[80]设计的 Hessian-free(HF)优化算法也可用于深度模型训练.

在训练过程中,深度学习模型用于控制模型结构和训练过程的超参数(hyper parameters)明显多于其他机器学习方法.重要的超参数包括隐藏层节点个数、激活函数选择、无监督或有监督过程、学习速率、训练次数、稀疏系数等.超参数的配置直接影响深度模型的训练效果,因此需要找到适当的超参数配置来进行训练才能得到最好的结果,目前这一问题主要依靠经验.但研究者也在积极寻找解决方案,如 Snoek 等^[81]提出的超参数自动优化方法,可根据当前训练结果和训练历史自动在超参数空间中搜寻最优超参数配置,避免了繁琐且盲目的手动调节过程.

如果在输入数据中加入随机噪声进行训练^[82],可以得到更加鲁棒的模型,这一特性非常适合数据噪声和缺失较多的生物医学问题.此外,深度学习在小数据集上容易出现过拟合问题.对此,一种有效的方法是 Hinton 等^[83]提出的 Dropout 技术,通过在训练过程中随机剔除神经元并查看结果,能够有效避免过拟合问题且能达到集成训练的效果.深度学习模型复杂且参数规模庞大,实际训练中比较耗时,因此最好使用图形处理器(GPU)或大规模并行系统进行并行加速.

3.3 深度学习应用的步骤

若采用深度学习技术解决生物医学大数据分析问题,首先应进行数据预处理,由于深度学习模型对输入数据的预设要求更低,因此预处理过程相比传统模型要简单,随后构建深度模型,针对具体问题使用不同的模型扩展及组合方法.最后训练深度模型,从原始数据中提取更高级的数据特征,提取到数据特征之后,可直接将其输入分类器进行分类,也可将其和原始信息一同作用于分类器.已有的实验表明,将处理后特征信息和原始数据信息共

同作为模型输入将得到更好的表现[34].

此外,值得一提的是图像处理中的"迁移学习",即使用通用图像数据训练模型,并将它用于另一图像集的处理. 如有研究者直接使用在ImageNet 数据集上训练好的 OverFeat 模型^[84],从小鼠脑部发育图片中提取特征^[61],效果很好. 这反映了深度学习模型和人类视觉系统在图像处理中的共性,"迁移学习"可以解决某些生物学问题中图像数据不足的困难,同时有效减少了深度学习模型在应用中的时间花费.

3.4 深度学习工具的使用

目前已有很多深度学习开源框架可供使用,降低了科研人员使用深度学习技术的门槛,未来还会有更多的深度学习工具出现.表2列出了七款笔者认为最具代表性的深度学习开源框架.

DeepLearnToolbox [85]是 MATLAB 工具包,它实现了除 RNN 之外的所有常见深度学习模型,其代码简单,适合初学者学习. Caffe 由贾扬清等[86] 开发,是目前最快的 CNN 实现,被广泛应用在计算机视觉领域,但 Caffe 仅实现了 CNN 模型. Torch [87]是基于 Lua 语言的深度学习框架,其运行速度极快,被广泛使用,但 Torch 未提供其他语言接口,因此拓展性差. Theano [88]由蒙特利尔大学开发,本身是一个可用于 GPU 的对多维数组进行高

效运算的 Python 库,被广泛用于神经网络构建,很多深度学习 Python 框架都基于 Theano 开发,如 Keras^[89]、Lasagne、Pylearn2^[90]等. DeepLearning4j^[91]是一款用 JAVA 实现的商用深度学习框架,与 Hadoop 和 Spark 紧密结合,运行效率高且支持大规模分布式. TensorFlow^[92]是 Google 公司的第二代深度学习框架,实现了所有的主流深度学习模型,稳定且易用,但目前开源的部分运算效率不高,且不支持分布式系统. Keras^[89]是基于 Theano^[88]的二次开发框架,其设计高度模块化,因此编程效率高,Keras 实现了 CNN 和 RNN 模型,出现时间短但发展十分迅速. MXNet^[93]是最新出现的 cxxnet 的下一代开源框架,提供多种语言接口,且支持大规模分布式并行,MXNet 的编程效率和运行效率都相当高.

另外表 2 中列出的所有框架均支持 GPU 加速,可大大缩短模型训练所需时间.除 TensorFlow 外,所有模型都同时支持 Windows 和 Linux 平台.模型训练所需的硬件资源则视问题规模而定,小型深度网络可用个人电脑训练,而生物学和医学中常见的大规模问题则要采用计算机集群来应对.另外,表 2 所列仅是各框架的官方版本,实际应用中研究者可根据需要查找相关框架的改进版本,甚至自己实现新的深度学习框架.

Table 2	The most popular frameworks for deep lear		
	表 2	最受欢迎的深度学习框架	

5言接口	分布式支持	首发时间	实现模型
IATLAB	不支持	2011	SAE/DBN/CNN
Python	不支持	2011	无
thon / MATLAB	不支持	2013	CNN
Lua	不支持	2013	SAE/DBN/CNN/RNN
JAVA	支持	2013	SAE/DBN/CNN
+ / Python	不支持	2015	SAE/DBN/CNN/RNN
Python	不支持	2015	CNN/RNN
R / C++ / Julia	支持	2015	CNN/RNN
	/ C++ / Julia		

所有框架均支持 GPU 加速, TensorFlow 只支持 Linux, 其余框架都支持 Windows 和 Linux 两种平台.

4 展 望

在传统生物医学数据分析中,一般将各领域专家手工提取的特征作为模型输入,大数据情境中这种方法存在两个问题: a. 生物医学数据的特征集往往依赖数据集,例如基因表达数据的特征依赖具

体的组织器官,因此实际应用中要对不同数据集分别提取特征.b.人工提取特征依赖先验知识,很难提取到潜在的复杂特征,而这类特征往往对分类至关重要.数据特征提取已逐渐成为生物医学数据分析的瓶颈,而深度学习方法能够很好地应对.

此外,一方面生物医学中有大量医疗影像数

据,如组织切片影像、核磁共振影像、X 光影像等,包含大量潜在的人体健康数据.这部分数据目前普遍依赖专家读取和分析,易受干扰且效率不高,并且人工分析只能对其浅层挖掘,造成数据资源浪费.深度学习技术,尤其是 CNN 模型,在图片分类、关键部位发现等方面表现优异.经过通用图像训练得到的深度学习模型可稍加修改应用于生物医学图像问题,能够更加快速准确地进行疾病诊断,将成为强有力的生物医学图像分析工具.

深度学习方法借助其无监督学习过程及多层结构,能自动从复杂原始数据中提取抽象特征. 当同一问题中存在多种数据集时,深度学习能针对不同数据集提取不同特征,如果同时使用多个数据集,则能捕捉到问题中有效的通用特征. 总之,深度学习的自动特征提取具有优秀的快速泛化能力,节省特征提取成本的同时,提升了分类效果,为突破大数据分析的瓶颈提供了方法.

军事医学科学院伯晓晨课题组曾将深度学习应用于脑卒中预测及诊断问题,通过调查患者的生活习惯既往病史等,结合使用颈部血管筛查等自动化测试数据来预测其患有脑卒中的概率,自动分析技术与自动检测的对接为自动诊断提供了技术支撑,在未来大有可为. 另外,我们还将探索深度学习技术在大规模生物数据分析中的应用,如 lincs 数据分析中,多细胞系情境对分析精度有更高的要求,以前使用的机器学习无能为力,而深度学习技术则大有希望,基因组序列的分析中原始数据的维数过高,直接使用则非常耗时,若能采用深度学习技术先提取抽象特征,则有希望在提高准确率的同时节省时间.

与以往的机器学习相比,将深度学习应用于生物医学数据分析也存在一些问题: a. 模型不易分析,这是神经网络模型的通病. 即在得到效果优秀的模型之后,要分析不同数据特征在该模型中的重要程度并不容易. b. 要求大量数据,深度学习模型由于复杂的模型结构,数据量少时容易出现过拟合现象,而生物中由于实验费用高,能提供的有标签数据很少,为深度学习模型的训练增加了困难. c. 运算开销大,深度学习模型规模大,训练时对内存和 CPU 都有较高要求,且模型训练时间长,因此对硬件环境要求较高,最好有计算机集群提供计算资源.

深度学习方法的出现已有多年,但在生物医学领域的应用大多仍处在入门级,模型简单且层数较

少.这一方面是由生物医学分析技术上固有的滞后性所决定的,另一方面则是因为缺乏通用的深度学习框架.目前研究者使用最多的仍是 SAE 和DBN,但其他深度学习模型也具有相当大的潜力,如目前广泛用于时间序列处理的 RNN 模型很契合生物序列处理问题的特点,RNN 模型在生物医学数据分析中的潜力需要更多研究者的挖掘.另外,如果能有一种适用不同类型生物医学问题的深度学习平台或者框架,则能帮助研究者快速使用深度学习这项新技术,更好地应对已经到来的生物医学大数据时代.

参考文献

- [1] Genomes Project C, Auton A, Brooks L D, *et al*. A global reference for human genetic variation. Nature, 2015, **526**(7571): 68–74
- [2] Consortium E P. The ENCODE (ENCyclopedia Of DNA Elements)Project. Science, 2004, 306(5696): 636–640
- [3] Chadwick L H. The NIH roadmap epigenomics program data resource. Epigenomics, 2012, 4(3): 317–324
- [4] Duan Q, Flynn C, Niepel M, et al. LINCS Canvas Browser: interactive web app to query, browse and interrogate LINCS L1000 gene expression signatures. Nucleic Acids Research, 2014, 42(Web Server issue): W449–460
- [5] Barrett T, Wilhite S E, Ledoux P, et al. NCBI GEO: archive for functional genomics data sets—update. Nucleic Acids Research, 2013, 41(Database issue): D991–995
- [6] Tomczak K, Czerwinska P, Wiznerowicz M. The cancer genome atlas (TCGA): an immeasurable source of knowledge. Contemporary Oncology, 2015, 19(1A): A68-77
- [7] Qin J, Li Y, Cai Z, *et al.* A metagenome-wide association study of gut microbiota in type 2 diabetes. Nature, 2012, **490**(7418): 55–60
- [8] Mardis E R. The impact of next-generation sequencing technology on genetics. Trends in Genetics: TIG, 2008, **24**(3): 133–141
- [9] May M. Life science technologies: Big biological impacts from big data. Science, 2014, 344(6189): 1298–1300
- [10] Hinton G E, Salakhutdinov R R. Reducing the dimensionality of data with neural networks. Science (New York, NY), 2006, 313(5786): 504–507
- [11] Farabet C, Couprie C, Najman L, et al. Learning hierarchical features for scene labeling. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2013, 35(8): 1915–1929
- [12] Krizhevsky A, Sutskever I, Hinton G E. Imagenet classification with deep convolutional neural networks. Advances in Neural Information Processing Systems, 2012: 1097–1105
- [13] Hinton G, Deng L, Yu D, et al. Deep neural networks for acoustic modeling in speech recognition: the shared views of four research groups. IEEE Signal Processing Magazine, 2012, 29(6): 82–97
- [14] Mikolov T, Deoras A, Povey D, et al. Strategies for training large scale neural network language models. 2011 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU), 2011:

- 196-201
- [15] Felleman D J, Essen D C V. Distributed hierarchical processing in the primate cerebral cortex. Cerebral Cortex, 1991, 1(1): 1–47
- [16] Deng L. Deep learning: methods and applications. Foundations and Trends in Signal Processing, 2014, **7**(3–4): 197–387
- [17] Hinton G E, Osindero S, Teh Y W. A fast learning algorithm for deep belief nets. Neural Computation, 2006, **18**(7): 1527–1554
- [18] Bengio Y, Lamblin P, Popovici D, et al. Greedy layer-wise training of deep networks//Schölkopf B, Platt J, Hoffman T. Advances in Neural Information Processing Systems. USA: MZT Press, 2007: 153–160
- [19] Lecun Y, Bottou L, Bengio Y, et al. Gradient-based learning applied to document recognition. Proceedings of the IEEE, 1998, 86(11): 2278–2324
- [20] Hihi S E, Bengio Y. Hierarchical Recurrent Neural Networks for Long-Term Dependencies, 1996
- [21] Hochreiter S, Schmidhuber J. Long short-term memory. Neural Computation, 1997, **9**(8): 1735–1780
- [22] Pascanu R, Mikolov T, Bengio Y. On the difficulty of training Recurrent Neural Networks. ArXiv e-prints, 2012, 1211: 5063
- [23] Sutskever I, Martens J, Hinton G. Generating Text with Recurrent Neural Networks, 2011: 1017–1024
- [24] Bengio Y, Laufer E, Alain G, *et al.* Deep Generative Stochastic Networks Trainable by Backprop, 2014: 226–234
- [25] Le Q V, Zou W Y, Yeung S Y, et al. Learning Hierarchical Invariant Spatio-temporal Features for Action Recognition with Independent Subspace Analysis. Proceedings of the 2011 IEEE Conference on Computer Vision and Pattern Recognition. Washington, DC, USA; IEEE Computer Society, 2011: 3361–3368 %@ 3978-3361-4577-0394-3362
- [26] Poon H, Domingos P. Sum-product networks: A new deep architecture. 2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops), 2011: 689-690
- [27] Learned-Miller E, Lee H, Huang G B. Learning hierarchical representations for face verification with convolutional deep belief networks. 2014 IEEE Conference on Computer Vision and Pattern Recognition. Los Alamitos, CA, USA; IEEE Computer Society. 2012: 2518–2525
- [28] Lee H, Pham P, Largman Y, et al. Unsupervised feature learning for audio classification using convolutional deep belief networks. 2009, 1096–1104
- [29] Vincent P, Larochelle H, Bengio Y, et al. Extracting and Composing Robust Features with Denoising Autoencoders. Proceedings of the 25th International Conference on Machine Learning. New York, NY, USA; ACM. 2008: 1096–1103 % @ 1978-1091-60558-60205-60554
- [30] Szegedy C, Liu W, Jia Y, *et al*. Going Deeper with Convolutions. arXiv:14094842 [cs], 2014
- [31] Chilimbi T, Suzue Y, Apacible J, et al. Project Adam: Building an Efficient and Scalable Deep Learning Training System. Proceedings of the 11th USENIX Conference on Operating Systems Design and Implementation. Berkeley, CA, USA; USENIX Association, 2014:

- 571-582
- [32] Hubel D H, Wiesel T N. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. The Journal of Physiology, 1962, **160**(1): 106–154
- [33] Wulsin D F, Gupta J R, Mani R, et al. Modeling electroencephalography waveforms with semi-supervised deep belief nets: fast classification and anomaly measurement. Journal of Neural Engineering, 2011, 8(3): 036015
- [34] Chakdar K, Potetz B. Deep Learning for the Semiautomated Analysis of Pap Smears. Medical Applications of Artificial Intelligence. CRC Press, 2014: 193–213
- [35] Gao X, Lin S, Wong T Y. Automatic feature learning to grade nuclear cataracts based on deep learning. IEEE Transactions on Biomedical Engineering, 2015, **62**(11): 2693–2701
- [36] Yang Z, Zhong S, Carass A, et al. Deep Learning for Cerebellar Ataxia Classification and Functional Score Regression. Machine learning in medical imaging MLMI (Workshop), author, 2014, 8679, 68–76
- [37] Ithapu V K, Singh V, Okonkwo O C, et al. Imaging-based enrichment criteria using deep learning algorithms for efficient clinical trials in mild cognitive impairment. Alzheimer's & Dementia, 2015, 11(12): 1489–1499
- [38] Suk H-I, Lee S-W, Shen D, *et al.* Hierarchical feature representation and multimodal fusion with deep learning for AD/MCI diagnosis. NeuroImage, 2014, **101**: 569–582
- [39] Suk H-I, Lee S-W, Shen D, *et al.* Deep sparse multi-task learning for feature selection in Alzheimer's disease diagnosis. Brain Structure and Function, 2015: 1–19
- [40] Li F, Tran L, Thung K-H, *et al*. A robust deep model for improved classification of AD/MCI patients. IEEE Journal of Biomedical and Health Informatics, 2015, **19**(5): 1610–1616
- [41] Suk H-I, Shen D. Deep learning-based feature representation for AD/MCI classification. Medical image computing and computer-assisted intervention: MICCAI International Conference on Medical Image Computing and Computer-Assisted Intervention, 2013, 16(Pt 2): 583–590
- [42] Ithapu V K, Singh V, Okonkwo O C, et al. Imaging-based enrichment criteria using deep learning algorithms for efficient clinical trials in mild cognitive impairment. Alzheimer's & Dementia: The Journal of the Alzheimer's Association, 2015, 11(12): 1489–1499
- [43] Suk H-I, Lee S-W, Shen D, *et al.* Deep sparse multi-task learning for feature selection in Alzheimer's disease diagnosis. Brain Structure & Function, 2015
- [44] Suk H-I, Lee S-W, Shen D, et al. Latent feature representation with stacked auto-encoder for AD/MCI diagnosis. Brain Structure & Function, 2015, 220(2): 841–859
- [45] Fakoor R, Ladhak F. Using deep learning to enhance cancer diagnosis and classification. Atlanta, Georgia, USA, 2013
- [46] Cruz-Roa A A, Arevalo Ovalle J E, Madabhushi A, *et al.* A deep learning architecture for image representation, visual interpretability and automated basal-cell carcinoma cancer

- detection. Medical image computing and computer-assisted intervention: MICCAI International Conference on Medical Image Computing and Computer-Assisted Intervention, 2013, **16** (Pt 2): 403–410
- [47] Xu J, Xiang L, Hang R, et al. Stacked Sparse Autoencoder (SSAE) based framework for nuclei patch classification on breast cancer histopathology. 2014 IEEE 11th International Symposium on Biomedical Imaging (ISBI). 2014: 999–1002
- [48] Li R, Zhang W, Suk H-I, et al. Deep learning based imaging data completion for improved brain disease diagnosis. Medical image computing and computer-assisted intervention: MICCAI International Conference on Medical Image Computing and Computer-Assisted Intervention, 2014, 17(Pt 3): 305–312
- [49] Gao X, Lin S, Wong T. Automatic Feature Learning to Grade Nuclear Cataracts Based on Deep Learning. IEEE transactions on bio-medical engineering, 2015
- [50] Liu G-P, Yan J-J, Wang Y-Q, et al. Deep learning based syndrome diagnosis of chronic gastritis. Computational and Mathematical Methods in Medicine, 2014(2014): 938350
- [51] Liao S, Gao Y, Oto A, et al. Representation learning: a unified deep learning framework for automatic prostate MR segmentation. Medical image computing and computer-assisted intervention: MICCAI International Conference on Medical Image Computing and Computer-Assisted Intervention, 2013, 16(Pt 2): 254–261
- [52] Prasoon A, Petersen K, Igel C, et al. Deep feature learning for knee cartilage segmentation using a triplanar convolutional neural network. Medical image computing and computer-assisted intervention: MICCAI International Conference on Medical Image Computing and Computer-Assisted Intervention, 2013, 16 (Pt 2): 246–253
- [53] Song Y, Zhang L, Chen S, et al. A deep learning based framework for accurate segmentation of cervical cytoplasm and nuclei//Engineering in Medicine and Biology Society (EMBC). 2014 36th Annual International Conference of the IEEE. USA: IEEE, 2014: 2903–2906
- [54] Carneiro G, Nascimento J C. Combining multiple dynamic models and deep learning architectures for tracking the left ventricle endocardium in ultrasound data. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2013, 35(11): 2592–2607
- [55] Cire?an D C, Giusti A, Gambardella L M, et al. Mitosis detection in breast cancer histology images with deep neural networks. Medical image computing and computer-assisted intervention: MICCAI International Conference on Medical Image Computing and Computer-Assisted Intervention, 2013, 16(Pt 2): 411–418
- [56] Majumdar A. Real-time Dynamic MRI Reconstruction using Stacked Denoising Autoencoder. 2015
- [57] Brosch T, Tam R, Initiative for the Alzheimers Disease N. Manifold learning of brain MRIs by deep learning. Medical image computing and computer-assisted intervention: MICCAI International Conference on Medical Image Computing and Computer-Assisted Intervention, 2013, 16(Pt 2): 633-640
- [58] Brosch T, Yoo Y, Li D K B, et al. Modeling the variability in brain

- morphology and lesion distribution in multiple sclerosis by deep learning. Medical image computing and computer-assisted intervention: MICCAI International Conference on Medical Image Computing and Computer-Assisted Intervention, 2014, **17** (Pt 2): 462–469
- [59] Zhang W, Li R, Deng H, et al. Deep convolutional neural networks for multi-modality isointense infant brain image segmentation. NeuroImage, 2015, 108: 214–224
- [60] Wu G, Kim M, Wang Q, et al. Unsupervised deep feature learning for deformable registration of MR brain images. Medical image computing and computer-assisted intervention: MICCAI International Conference on Medical Image Computing and Computer-Assisted Intervention, 2013, 16(Pt 2): 649–656
- [61] Zeng T, Li R, Mukkamala R, *et al.* Deep convolutional neural networks for annotating gene expression patterns in the mouse brain. BMC Bioinformatics, 2015, **16**(1): 147–156
- [62] Helmstaedter M, Briggman K L, Turaga S C, et al. Connectomic reconstruction of the inner plexiform layer in the mouse retina. Nature, 2013, 500(7461): 168–174
- [63] Wang D, Shang Y. Modeling physiological data with deep belief networks. International Journal of Information and Education Technology (IJIET), 2013, 3(5): 505-511
- [64] Lyons J, Dehzangi A, Heffernan R, et al. Predicting backbone Cα angles and dihedrals from protein sequences by stacked sparse auto-encoder deep neural network. Journal of Computational Chemistry, 2014, 35(28): 2040–2046
- [65] Spencer M, Eickholt J, Cheng J. A deep learning network approach to ab initio protein secondary structure prediction. IEEE/ACM Transactions on Computational Biology and Bioinformatics/IEEE, ACM, 2015, 12(1): 103–112
- [66] Heffernan R, Paliwal K, Lyons J, et al. Improving prediction of secondary structure, local backbone angles, and solvent accessible surface area of proteins by iterative deep learning. Scientific Reports, 2015, 5: 11476
- [67] Di Lena P, Nagata K, Baldi P. Deep architectures for protein contact map prediction. Bioinformatics (Oxford, England), 2012, 28 (19): 2449–2457
- [68] Lena P D, Nagata K, Baldi P F. Deep Spatio-Temporal Architectures and Learning for Protein Structure Prediction [M]. Harrahs and Harveys, Lake Tahoe. 2012
- [69] Zhou J, Troyanskaya O. Deep Supervised and Convolutional Generative Stochastic Network for Protein Secondary Structure Prediction. 2014: 745–753
- [70] Eickholt J, Cheng J. Predicting protein residue-residue contacts using deep networks and boosting. Bioinformatics (Oxford, England), 2012, 28(23): 3066–3072
- [71] Jo T, Hou J, Eickholt J, *et al.* Improving protein fold recognition by deep learning networks. Scientific Reports, 2015, **5**: 17573
- [72] Quang D, Chen Y, Xie X. DANN: a deep learning approach for annotating the pathogenicity of genetic variants. Bioinformatics (Oxford, England), 2015, **31**(5): 761–763
- [73] Leung M K K, Xiong H Y, Lee L J, et al. Deep learning of the

- tissue-regulated splicing code. Bioinformatics (Oxford, England), 2014, **30**(12): i121–129
- [74] Xiong H Y, Alipanahi B, Lee L J, et al. RNA splicing. The human splicing code reveals new insights into the genetic determinants of disease. Science, 2015, 347(6218): 1254806
- [75] Fan X N, Zhang S W. IncRNA-MFDL: identification of human long non-coding RNAs by fusing multiple features and using deep learning. Molecular Biosystems, 2015, 11(3): 892–897
- [76] Liang M, Li Z, Chen T, et al. Integrative data analysis of multi-platform cancer data with a multimodal deep learning approach. IEEE/ACM Transactions on Computational Biology and Bioinformatics/IEEE, ACM, 2015, 12(4): 928–937
- [77] Ibrahim R, Yousri N A, Ismail M A, et al. Multi-level gene/MiRNA feature selection using deep belief nets and active learning. Conference proceedings: Annual International Conference of the IEEE Engineering in Medicine and Biology Society IEEE Engineering in Medicine and Biology Society Annual Conference, 2014, 2014: 3957–3960
- [78] Ma J, Sheridan R P, Liaw A, et al. Deep neural nets as a method for quantitative structure-activity relationships. Journal of Chemical Information and Modeling, 2015, 55(2): 263–274
- [79] Nguyen S P, Shang Y, Xu D. DL-PRO: A Novel Deep Learning Method for Protein Model Quality Assessment. Proceedings of International Joint Conference on Neural Networks / co-sponsored by Japanese Neural Network Society (JNNS) [et al.] International Joint Conference on Neural Networks, 2014, 2014: 2071–2078
- [80] Martens J. Deep learning via hessian-free optimization. Haifa, Israel. 2010: 735–742
- [81] Snoek J, Larochelle H, Adams R P. Practical Bayesian optimization of machine learning algorithms. In Advances in Neural Information

- Processing Systems. 2012
- [82] Bengio Y, Yao L, Alain G, *et al.* Generalized Denoising Auto-Encoders as Generative Models. arXiv:13056663 [cs], 2013
- [83] Srivastava N, Srivastava N. Improving Neural Networks with Dropout. University of Toronto, 2013
- [84] Sermanet P, Eigen D, Zhang X, et al. OverFeat: Integrated Recognition, Localization and Detection using Convolutional Networks. arXiv:13126229 [cs], 2013
- [85] Palm R B. Prediction as a Candidate for Learning Deep Hierarchical Models of Data. Technical University of Denmark, 2012
- [86] Jia Y, Shelhamer E, Donahue J, et al. Caffe: Convolutional Architecture for Fast Feature Embedding. Eprint Arxiv, 2014, 675–678
- [87] Collobert R, Bengio S, Marithoz J. Torch: A Modular Machine Learning Software Library. Idiap, 2002
- [88] Bastien F, Lamblin P, Pascanu R, *et al.* Theano: new features and speed improvements. Eprint Arxiv, 2012
- [89] Keras Documentation, http://keras.io
- [90] Goodfellow I J, Warde-Farley D, Lamblin P, *et al.* Pylearn2: a machine learning research library. Eprint Arxiv, 2013
- [91] Team D J D. Deeplearning4j: Open-source distributed deep learning for the JVM, Apache Software Foundation License 2.0
- [92] Martín Abadi A A, Paul Barham, Eugene Brevdo, et al. TensorFlow: Large-scale machine learning on heterogeneous systems. 2015
- [93] Tianqi Chen M L, Yutian Li, Min Lin, et al. MXNet: A Flexible and Efficient Machine Learning Library for Heterogeneous Distributed Systems. In Neural Information Processing Systems, Workshop on Machine Learning Systems, 2015

Applications of Deep Learning in Biological and Medical Data Analysis*

LI Yuan¹, LUO Zhi-Gang¹, GUAN Nai-Yang¹, YIN Xiao-Yao¹, WANG Bing³, BO Xiao-Chen², LI Fei²,

(1) Science and Technology on Parallel and Distributed Processing Laboratory, College of Computer,
National University of Defense Technology, Changsha 410073, China;

2) Institute of Radiation Medicine, Academy of Military Medical Sciences, Beijing 100850, China;

3) Unit 63928 of Chinese People's Liberation Army, Beijing 100028, China)

Abstract The rapid accumulation of biomedical data provided unprecedented opportunities for biology and clinical research, while it also made traditional data analysis technology face enormous challenges. In this paper, we reviewed recent studies on biomedical data using deep learning. We introduced several recommended deep learning models and summarized current applications of biological and medical data analysis using deep learning, including the general procedure, model construction and training process. Finally, we made a discussion on some issues in deep learning applications.

Key words deep learning, high-throughput omics, clinical medicine, data mining **DOI**: 10.16476/j.pibb.2015.0339

LUO Zhi-Gang. Tel: 86-731-84575835, E-mail: zgluo@nudt.edu.cn

LI Fei. Tel: 86-10-66932251, E-mail: pittacus@gmail.com Received: January 21, 2016 Accepted: March 25, 2016

^{*} This work was supported by grants from The National High-tech R&D Program of China(2015AA020100), The National Natural Science Foundation of China(U1435222, 81273488, 61402486), Research in Key Technologies of Infectious Disease Prevention of China(BWS14C051).

^{**}Corresponding author.