



南京航空航天大学
NANJING UNIVERSITY OF AERONAUTICS AND ASTRONAUTICS

第 2 章 深度学习基础

魏明强、宫丽娜

计算机科学与技术学院

智周万物·道济天下

目录



- 神经网络基础
 - 神经网络
 - 卷积神经网络
- 损失函数和优化算法
 - 损失函数
 - 优化算法
- 神经网络训练
 - 梯度和链式法则
 - 前向传播和反向传播
 - 训练神经网络示例
- 深度学习框架
 - 主流深度学习框架
 - 框架选择和优缺点比较
- 思考



- 神经网络基础
 - 神经网络
 - 卷积神经网络
- 损失函数和优化算法
 - 损失函数
 - 优化算法
- 神经网络训练
 - 梯度和链式法则
 - 前向传播和反向传播
 - 训练神经网络示例
- 深度学习框架
 - 主流深度学习框架
 - 框架选择和优缺点比较
- 思考

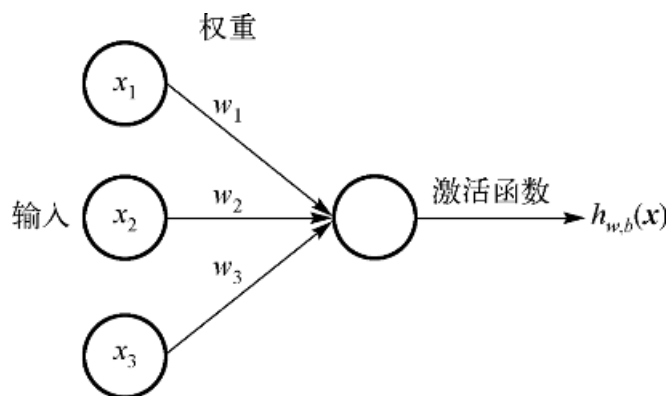
神经网络

□ 人工神经元

人工神经网络 (Artificial Neural Network, ANN), 简称为神经网络 (Neural Network: NN), 是指一系列受生物学和神经科学启发的数学模型。

人工神经元, 简称为神经元, 是构成神经网络的基本单元。

$$h_{w,b}(x) = f(w^T x + b)$$

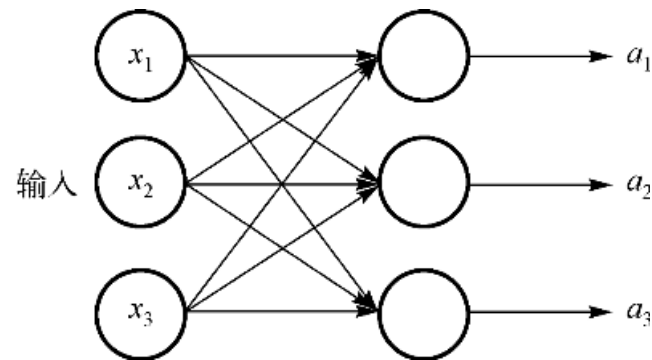


单个神经元计算过程

要想模拟人脑具有的能力, 单一神经元是远远不够的, 需要众多神经元的协作来完成复杂任务, 即神经网络。

在得到单层神经网络的输出之后, 可以通过叠加类似的层来构建每层都包含若干神经元的多层神经网络。

$$a = f(Wx + b)$$



单层神经网络计算过程

神经网络

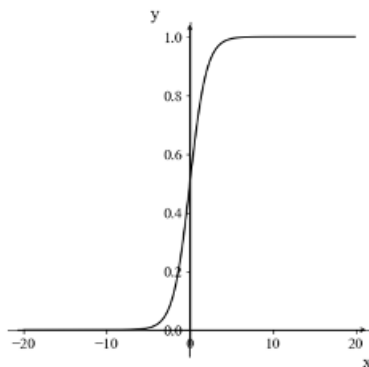
□ 激活函数

激活函数 (Activation Function) 是神经网络中的一种非线性变换, 它赋予神经元更强大的表达能力。如果不使用激活函数, 则每层的操作只是对上一层的输出结果进行线性变换, 多层神经网络会退化成单层神经网络。

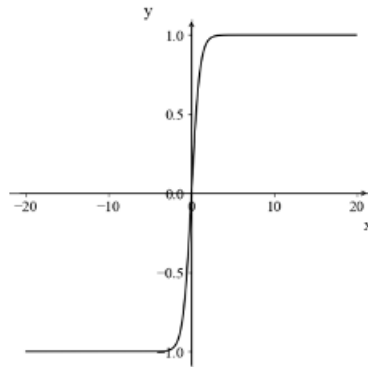
- Sigmoid函数

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

通常用于二分类问题的输出层。



(a) Sigmoid函数



(b) Tanh函数

- Tanh函数

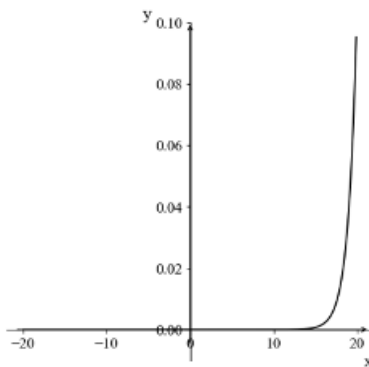
$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

通常用于中间层或输出层。

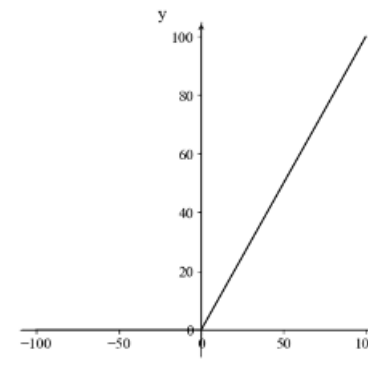
- Softmax函数

$$\text{softmax}(x)_i = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

通常用于多分类问题的输出层。



(c) Softmax函数



(d) ReLU函数

- ReLU函数

$$f_{\text{ReLU}}(x) = \max(0, x)$$

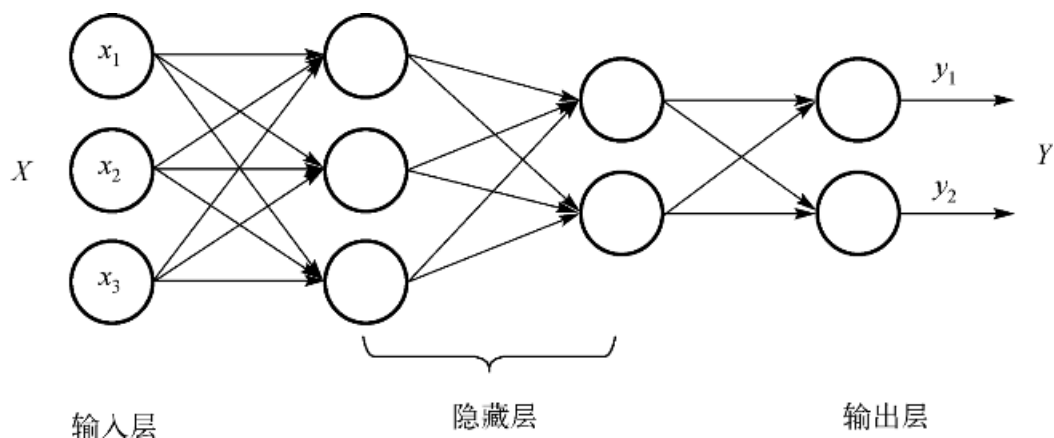
广泛应用于隐藏层, 其简单性和非饱和性使其在大多数情况下表现良好。

神经网络

□ 全连接神经网络

在全连接神经网络中，每个神经元与前一层的所有神经元相连接，形成一个完全连接的结构。

它的基本组成包括输入层（Input Layer）、若干隐藏层（Hidden Layer）和输出层（Output Layer）。输入层接收原始数据或特征作为网络的输入，每个输入神经元对应于数据或特征的一个维度。隐藏层位于输入层和输出层之间，进行特征的非线性变换和抽象。每个隐藏层包含多个神经元，每个神经元与前一层的所有神经元相连接。多个隐藏层的存在使得网络能够学习更加复杂和抽象的表示。输出层产生网络的最终输出。



全连接神经网络在一些任务上表现良好，但随着问题复杂性的增加，更深层次、更复杂结构的神经网络逐渐取代了全连接神经网络。这是因为全连接神经网络在参数数量和计算复杂度上容易受到限制，而深度学习任务通常需要更强大的神经网络结构。

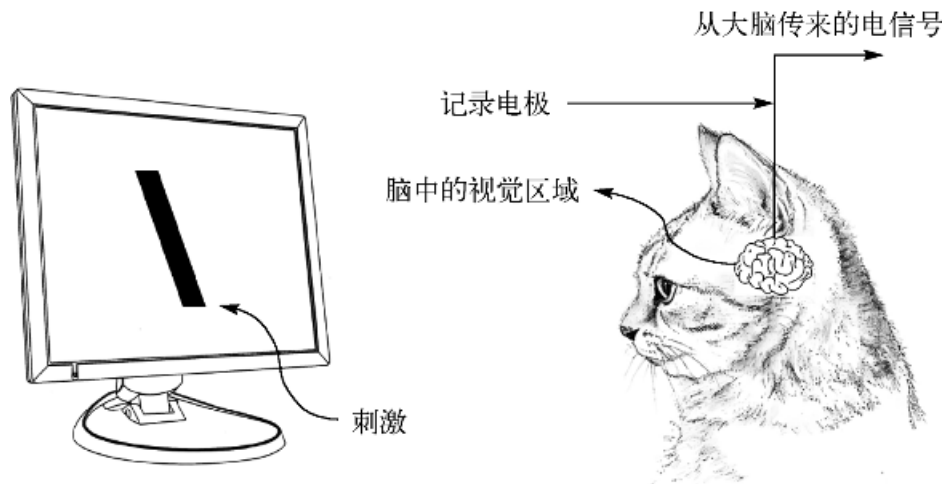


- 神经网络基础
 - 神经网络
 - 卷积神经网络
- 损失函数和优化算法
 - 损失函数
 - 优化算法
- 神经网络训练
 - 梯度和链式法则
 - 前向传播和反向传播
 - 训练神经网络示例
- 深度学习框架
 - 主流深度学习框架
 - 框架选择和优缺点比较
- 思考

卷积神经网络

□ 感受野

1962年，生物学家D.H. Hubel和T.N. Wiesel对猫的视觉系统进行了研究，猫的视觉系统实验示意图如图2.5所示。他们首次发现了在猫的视觉皮层中存在两种主要类型的神经元，即简单细胞和复杂细胞。这两种类型的细胞对边缘和纹理的敏感性有所不同。神经元对视野中的某一小块区域内的特定边缘或纹理更为敏感，反映了感受野的特性。



感受野 (Receptive Field) 描述了神经系统中一些神经元对于特定刺激区域的敏感性，这意味着神经元只对其支配区域内的信号做出响应。在视觉神经系统中，视觉皮层中的神经细胞的输出受到视网膜上光感受器的影响，即当视网膜上的光感受器受到刺激并兴奋时，会产生神经冲动信号并传递到视觉皮层。然而，并非所有视觉皮层中的神经元都会接收这些信号。每个神经元都有其特定的感受野，即只有视网膜上特定区域内的刺激才能激活该神经元。

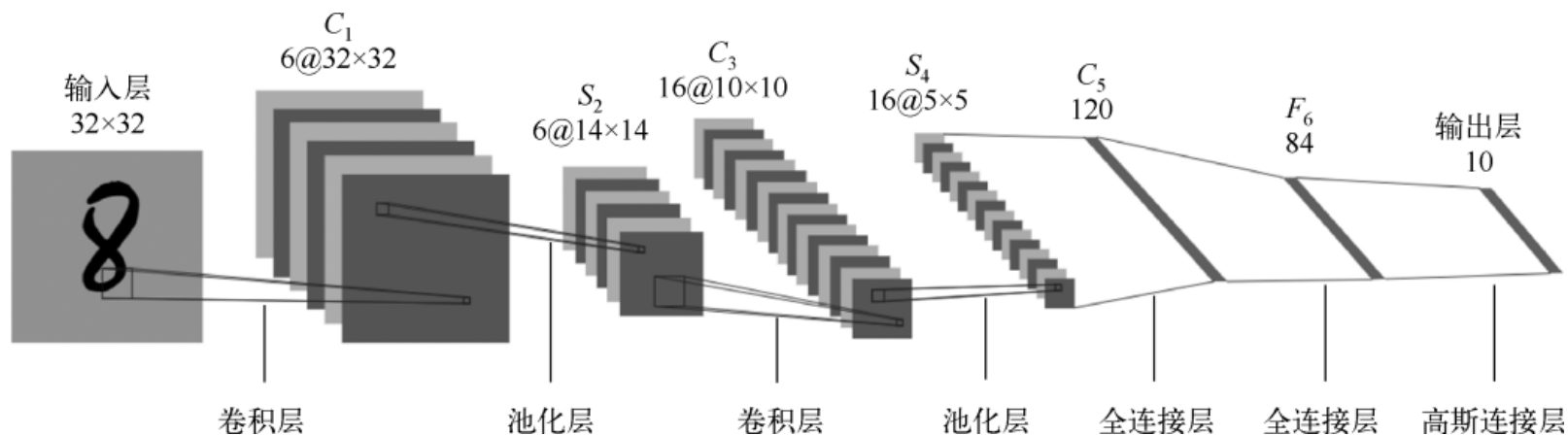
卷积神经网络

□ 卷积神经网络

卷积神经网络 (Convolutional Neural Network, CNN) 的设计灵感正是源自生物学中感受野的机制。

卷积神经网络模仿了生物学中神经元对于刺激的局部敏感性。它通过学习局部特征，逐渐建立对整体特征的抽象。它在处理空间结构化数据和视觉数据方面的能力使其在自然语言处理、计算机视觉等领域都发挥着重要作用。

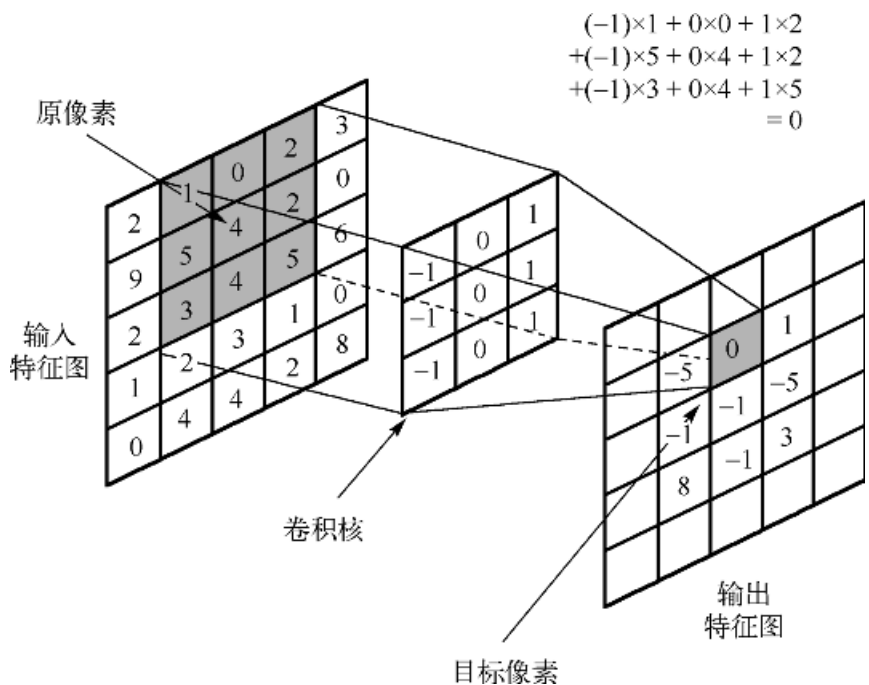
下图展示了第一个诞生的卷积神经网络LeNet-5的网络结构，该网络用于手写数字识别任务。LeNet-5由卷积层、池化层及全连接层组成，它的设计为后续卷积神经网络的发展奠定了基础。



卷积神经网络

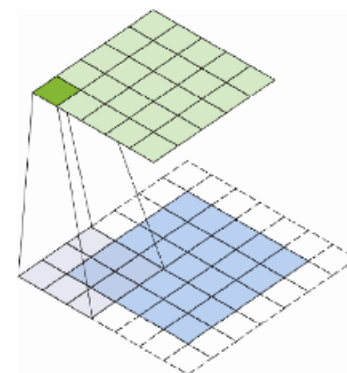
□ 卷积

卷积运算通过滑动一定间隔的卷积核（也称为滤波器）窗口，计算对应位置的元素相乘再求和，得到输出特征图中每个位置的值，当卷积核窗口移动到所示位置时，计算输入特征图与卷积核窗口对应位置的元素乘积，并将其求和，即执行计算： $(-1) \times 1 + 0 \times 0 + 1 \times 2 + (-1) \times 5 + 0 \times 4 + 1 \times 2 + (-1) \times 3 + 0 \times 4 + 1 \times 5 = 0$ ，从而计算得到输出特征图中相应位置的值为0。之后，卷积核继续向后滑动，重复相同的操作，直到得到完整的输出特征图。

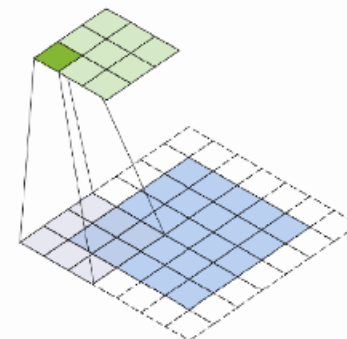


□ 卷积操作的概念

- 偏置 (bias)
- 步长 (stride)
- 填充 (padding)



(a) 步长=1



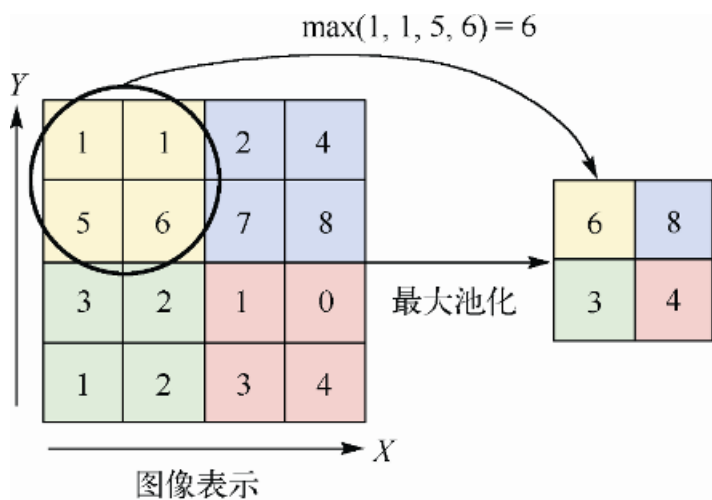
(b) 步长=2

$$\text{output size} = \frac{\text{input size} - \text{kernel size} + 2 \times \text{padding}}{\text{stride}} + 1$$

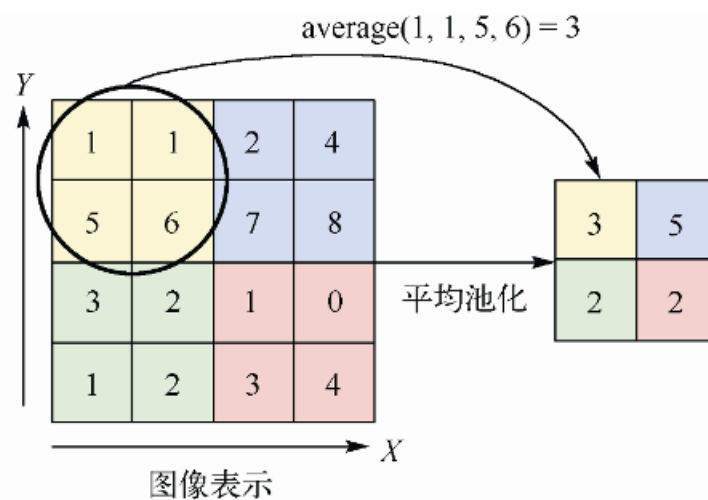
卷积神经网络

池化

池化操作通常应用在卷积层之后，通过对特征图的局部区域进行采样，从而获得更小且具有抽象特征的特征图。常见的池化类型有最大池化和平均池化两种。在最大池化中，每个池化窗口选择局部区域的最大值作为采样值。而在平均池化中，每个池化窗口计算局部区域的平均值作为采样值。



(a) 最大池化



(b) 平均池化

池化层的特点

- 没有可学习参数
- 不改变通道数
- 平移不变性

$$\text{output size} = \frac{\text{input size} - \text{padding kernel size}}{\text{stride}} + 1$$

卷积神经网络

□ 批归一化

批归一化的作用是加速神经网络的训练，提高模型的收敛速度，并且有助于避免梯度消失或梯度爆炸问题。批归一化的核心思想是对每层的输入进行归一化，使其均值接近 0，标准差接近 1。这样做有助于缓解梯度消失问题，提高网络的稳定性。对于一个批次的输入数据，批归一化首先计算批次的均值和方差，再对输入进行归一化，即减去均值并除以标准差，然后使用可学习的缩放和平移参数对归一化后的数据进行线性变换。

$$\text{BN}(x) = \gamma \frac{x - \mu}{\sigma} + \beta$$

□ 全连接

全连接层（Fully Connected Layer），也被称为密集连接层，是卷积神经网络中的关键组成部分。在全连接层中，每个神经元都与上一层的所有神经元相连接，形成了一个全连接的结构。对于自然语言处理任务，输入通常是一维向量，如文本数据的词嵌入，以便进行文本分类、情感分析等任务；对于计算机视觉任务，输入通常是多维特征图，这些特征图可能通过卷积层或其他特征提取层从原始图像中提取而来。为了传递给全连接层，这些多维特征图通常需要被展平成一维向量，作为全连接层的输入，以便进行后续处理。

卷积神经网络

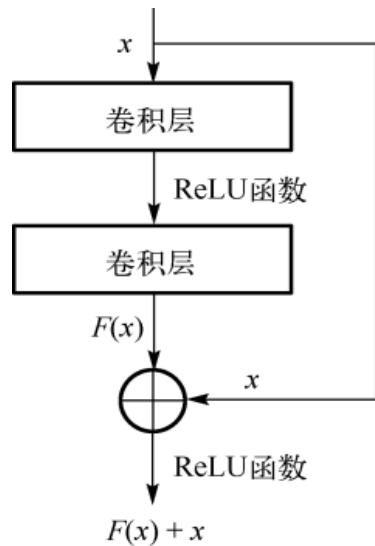
□ Dropout

Dropout是一种常用的正则化技术，旨在减少过拟合并提高模型的泛化能力。Dropout的基本思想是在训练过程中以一定概率随机地忽略一部分神经元的输出。具体而言，假设有一个全连接层的输出向量为 h ，Dropout的操作如下：

- (1) 在训练中，以概率（通常为 0.5）随机选择一部分神经元，将它们的输出置为0。
- (2) 在测试过程中，保持所有神经元的输出，但将它们乘以 $1 - p$ 以保持期望输出值不变。

□ 残差连接

残差连接将若干卷积层学习到的特征与原始输入相加，从而形成了一种“跳跃连接”的结构，从而使得神经网络更容易进行优化，并且能够构建更深层次的网络结构。残差连接能够在一定程度上缓解深层网络的退化网络问题。并且既不增加额外的参数也不增加计算复杂度，使得网络易于优化，提高了泛化性能。



目录



- 神经网络基础
 - 神经网络
 - 卷积神经网络
- 损失函数和优化算法
 - 损失函数
 - 优化算法
- 神经网络训练
 - 梯度和链式法则
 - 前向传播和反向传播
 - 训练神经网络示例
- 深度学习框架
 - 主流深度学习框架
 - 框架选择和优缺点比较
- 思考

损失函数

□ 均方误差损失函数

均方误差 (Mean Squared Error , MSE) 损失函数是一种应用于回归问题的损失函数 , 用于度量模型预测值与真实值之间的平方差的平均值。

$$\text{MSE}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

□ 平方绝对误差损失函数

平均绝对误差 (Mean Absolute Error , MAE) 损失函数是应用于回归问题的一种损失函数 , 用于度量模型预测值与真实值之间的绝对差的平均值。

$$\text{MAE}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

□ 交叉熵损失函数

交叉熵损失 (Cross-Entropy Loss) 函数广泛应用于分类问题。它衡量模型输出的概率分布与真实标签的概率分布之间的差异。

二分类问题 :

$$\text{Binary Cross-Entropy}(y, \hat{y}) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}))$$

多分类问题 :

$$\text{Categorical Cross-Entropy}(y, \hat{y}) = - \sum_{i=1}^C y_i \log(\hat{y}_i)$$

损失函数

□ 序列交叉熵损失函数

序列交叉熵损失 (Sequence Cross-Entropy Loss) 函数是用于序列到序列 (sequence- to-sequence) 任务中的一种损失函数，主要应用于自然语言处理领域的机器翻译任务。在这种任务中，模型需要将一个输入序列映射到另一个输出序列，而且输入和输出的序列长度是可变的

$$\text{Sequence Cross-Entropy Loss} = -\frac{1}{T} \sum_{t=1}^T \sum_{i=1}^N y_{t,i} \log(\hat{y}_{t,i})$$

□ 焦点损失函数

焦点损失 (Focal Loss) 函数通过调整难易分类样本的权重，即降低易分类样本的权重，提高难分类样本的权重，使得模型更关注难以分类的样本。

$$\text{Focal Loss} = -(1 - \hat{y}_t)^\gamma \log(\hat{y}_t)$$

$$\hat{y}_t = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{otherwise} \end{cases}$$

目录



- 神经网络基础
 - 神经网络
 - 卷积神经网络
- 损失函数和优化算法
 - 损失函数
 - **优化算法**
- 神经网络训练
 - 梯度和链式法则
 - 前向传播和反向传播
 - 训练神经网络示例
- 深度学习框架
 - 主流深度学习框架
 - 框架选择和优缺点比较
- 思考

优化算法

□ 梯度下降法

(1) 目标函数设定：设定一个目标函数（损失函数） $J(\theta)$ ，其中 θ 表示模型的参数。目标是找到使得目标函数最小化的参数值。

(2) 梯度计算：计算目标函数关于参数的梯度，即 $\nabla J(\theta)$ 。梯度表示目标函数在参数空间中的变化率，它指示了函数增长最快的方向。

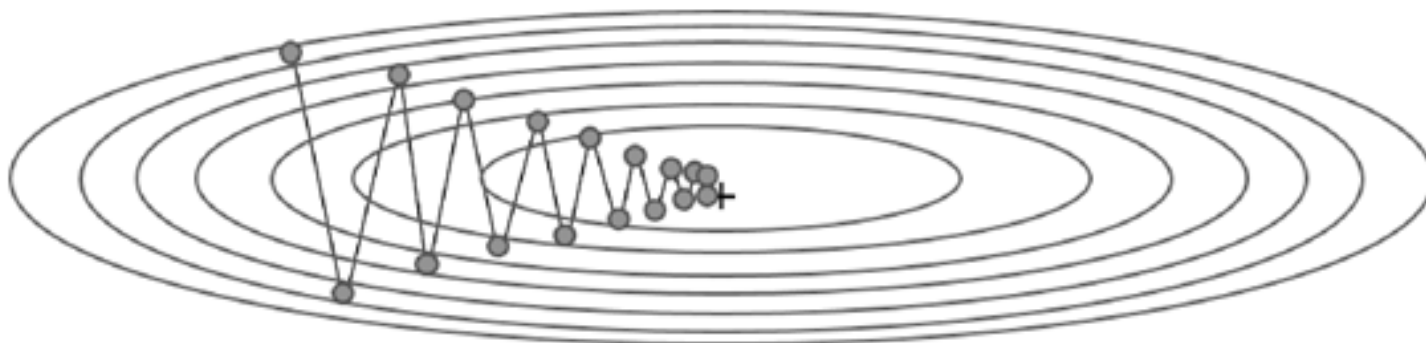
(3) 参数更新：根据梯度的反方向调整参数，通过以下规则进行参数更新，即

$$\theta = \theta - \alpha \nabla J(\theta)$$

其中， α 是学习率，表示每次更新时沿梯度方向移动的步长。学习率的选择对梯度下降的性能有很大的影响，学习率过大可能导致震荡或发散，学习率过小可能导致收敛缓慢。

□ 梯度下降法变种

- 批量梯度下降法
- 随机梯度下降法
- 小批量梯度下降法



优化算法

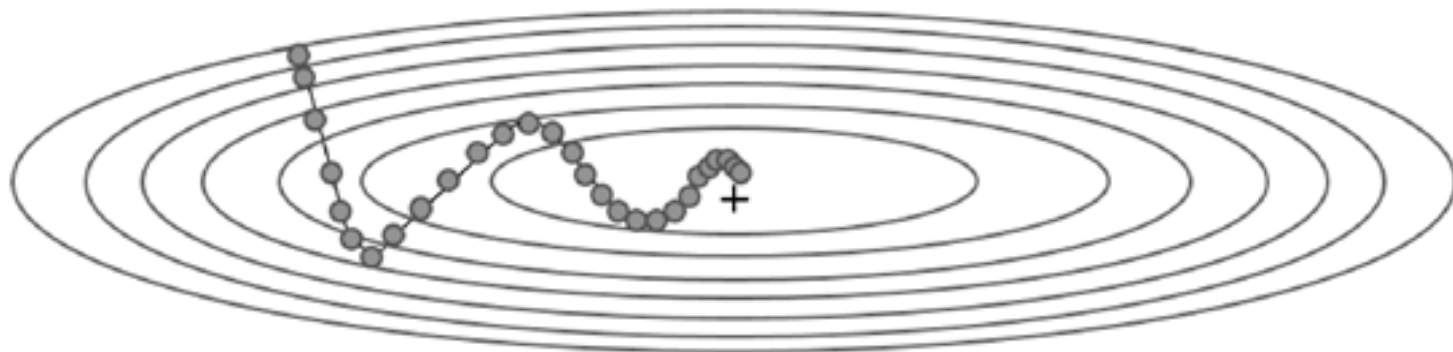
□ 动量法

动量法 (Momentum) 旨在加速收敛并减小震荡。它引入了梯度的“动量”或“速度”概念，类似于物理学中的动量。动量法的核心思想是在更新参数时，不仅考虑当前的梯度，还考虑过去梯度的累积效果，以减小在参数空间的震荡。

$$v = \beta v + (1 - \beta) \nabla J(\theta)$$

$$\theta = \theta - \alpha v$$

其中， v 是动量（速度）； β 是动量系数，控制过去梯度的权重，通常取值为(0, 1)，较大的 β 表示更多地考虑过去的梯度，从而减小震荡。通常 β 取值为0.9是一个常见的起点。 $\nabla J(\theta)$ 是目标函数关于参数的梯度。



以上内容仅为本文档的试下载部分，为可阅读页数的一半内容。如要下载或阅读全文，请访问：
<https://d.book118.com/728044117136007004>