

1.为什么引入非线性激励函数？

---

2.为什么引入relu

3.什么是好的激活函数？

## 1.为什么引入非线性激励函数？

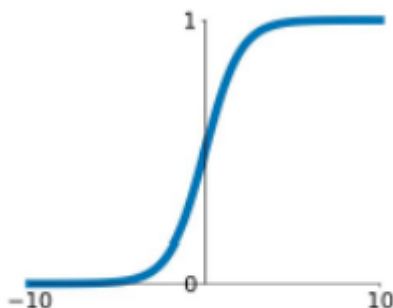
如果不用激励函数，在这种情况下每一层输出都是上层输入的线性函数，很容易验证，无论神经网络有多少层，输出都是输入的线性组合，与没有隐藏层效果相当，这种情况就是最原始的感知机（perception）了。

所以我们决定引入非线性函数作为激励函数，这样深层神经网络就有意义了（不再是输入的线性组合，可以逼近任意函数）。最早的想法是sigmoid或者tanh函数，输出有界，很容易充当下一层输入

## 2.为什么引入relu

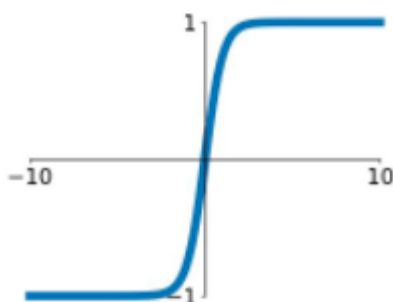
# Sigmoid

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



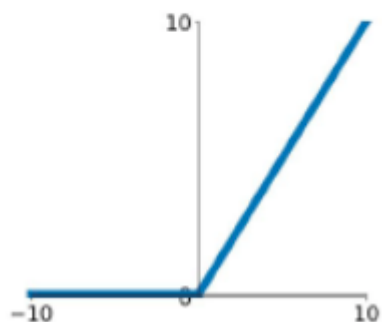
# tanh

$$\tanh(x)$$



# ReLU

$$\max(0, x)$$



第一，采用sigmoid等函数，算激活函数时（指数运算），计算量大，反向传播求误差梯度时，求导涉及除法，计算量相对大，而采用relu激活函数，整个过程的计算量节省很多。

第二，对于深层网络，sigmoid函数反向传播时，很容易出现梯度消失（在sigmoid接近饱和区时，变化太缓慢，导数趋于0），从而无法完成深层网络的训练。

第三，relu会使一部分神经元的输出为0，这样就造成了网络的稀疏性，并且减少了参数的相互依存关系，缓解了过拟合问题的发生。

当然现在也有一些对relu的改进，比如prelu, random relu等，在不同的数据集上会有一些训练速度上或者准确率上的改进，具体的大家可以找相关的paper看。

多加一句，现在主流的做法，会多做一步batch normalization，尽可能保证每一层网络的输入具有相同的分布

### 3.什么是好的激活函数？

作者：Hengkai Guo

链接：<https://www.zhihu.com/question/67366051/answer/262087707>

来源：知乎

著作权归作者所有。商业转载请联系作者获得授权，非商业转载请注明出处。

说说我对一个好的激活函数的理解吧，有些地方可能不太严谨，欢迎讨论。（部分参考了[Activation function](#)。）

1. **非线性**：即导数不是常数。这个条件前面很多答主都提到了，是多层神经网络的基础，保证多层网络不退化成单层线性网络。这也是激活函数的意义所在。
2. **几乎处处可微**：可微性保证了在优化中梯度的可计算性。传统的激活函数如sigmoid等满足处处可微。对于分段线性函数比如ReLU，只满足几乎处处可微（即仅在有限个点处不可微）。对于SGD算法来说，由于几乎不可能收敛到梯度接近零的位置，有限的不可微点对于优化结果不会有很大影响[1]。
3. **计算简单**：正如题主所说，非线性函数有很多。极端的说，一个多层神经网络也可以作为一个非线性函数，类似于Network In Network[2]中把它当做卷积操作的做法。但激活函数在神经网络前向的计算次数与神经元的个数成正比，因此简单的非线性函数自然更适合作为激活函数。这也是ReLU之流比其它使用Exp等操作的激活函数更受欢迎的其中一个原因。
4. **非饱和性（saturation）**：饱和指的是在某些区间梯度接近于零（即梯度消失），使得参数无法继续更新的问题。最经典的例子是Sigmoid，它的导数在x为比较大的正值和比较小的负值时都会接近于0。更极端的例子是阶跃函数，由于它在几乎所有位置的梯度都为0，因

此处饱和，无法作为激活函数。ReLU在 $x > 0$ 时导数恒为1，因此对于再大的正值也不会饱和。但同时对于 $x < 0$ ，其梯度恒为0，这时候它也会出现饱和的现象（在这种情况下通常称为dying ReLU）。Leaky ReLU[3]和PReLU[4]的提出正是为了解决这一问题。

**5. 单调性 (monotonic)：**即导数符号不变。这个性质大部分激活函数都有，除了诸如sin、cos等。个人理解，单调性使得在激活函数处的梯度方向不会经常改变，从而让训练更容易收敛。

**6. 输出范围有限：**有限的输出范围使得网络对于一些比较大的输入也会比较稳定，这也是为什么早期的激活函数都以此类函数为主，如Sigmoid、TanH。但这导致了前面提到的梯度消失问题，而且强行让每一层的输出限制到固定范围会限制其表达能力。因此现在这类函数仅用于某些需要特定输出范围的场合，比如概率输出（此时loss函数中的log操作能够抵消其梯度消失的影响[1]）、LSTM里的gate函数。

**7. 接近恒等变换 (identity)：**即约等于 $x$ 。这样的好处是使得输出的幅值不会随着深度的增加而发生显著的增加，从而使网络更为稳定，同时梯度也能够更容易地回传。这个与非线性是有点矛盾的，因此激活函数基本只是部分满足这个条件，比如TanH只在原点附近有线性区（在原点为0且在原点的导数为1），而ReLU只在 $x > 0$ 时为线性。这个性质也让初始化参数范围的推导更为简单[5][4]。额外提一句，这种恒等变换的性质也被其他一些网络结构设计所借鉴，比如CNN中的ResNet[6]和RNN中的LSTM。

**8. 参数少：**大部分激活函数都是没有参数的。像PReLU带单个参数会略微增加网络的大小。还有一个例外是Maxout[7]，尽管本身没有参数，但在同样输出通道数下k路Maxout需要的输入通道数是其它函数的k倍，这意味着神经元数目也需要变为k倍；但如果不考虑维持输出通道数的情况下，该激活函数又能将参数个数减少为原来的k倍。

9. 归一化 (normalization)：这个是最近才出来的概念，对应的激活函数是SELU[8]，主要思想是使样本分布自动归一化到零均值、单位方差的分布，从而稳定训练。在这之前，这种归一化的思想也被用于网络结构的设计，比如Batch Normalization[9]。