

第六讲 训练神经网络一

2019年4月30日 10:54

1. 激活函数

1.1 大纲

Overview

1. One time setup

activation functions, preprocessing, weight initialization, regularization, gradient checking

2. Training dynamics

babysitting the learning process, parameter updates, hyperparameter optimization

3. Evaluation

model ensembles

Stanf

图6.1.1 训练神经网络小节大纲

1. 搭建

激活函数、数据预处理、权重初始化、正则项、梯度检查

2. 训练中的动态变化

如何监测整个学习过程、参数(权重)更新、超参数优化

3. 模型评估

模型集成

1.2 part 1

Part 1

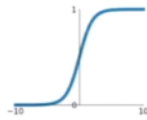
- Activation Functions
- Data Preprocessing
- Weight Initialization
- Batch Normalization
- Babysitting the Learning Process
- Hyperparameter Optimization

1.2.1 激活函数

Activation Functions

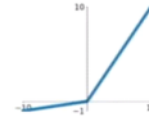
Sigmoid

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



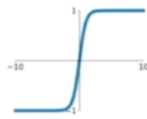
Leaky ReLU

$$\max(0.1x, x)$$



tanh

$$\tanh(x)$$

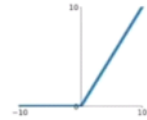


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ReLU

$$\max(0, x)$$



ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

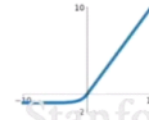


图6.1.2 常用激活函数

① sigmoid: historically popular!

- 缺点: 1.会使饱和神经元梯度消失
2.输出不是以0为中心的分布
3.函数中包含exp运算, 代价太大

当输入远大于10的时候sigmoid函数的梯度接近0

② tanh: 优点: 0中心分布

缺点: 存在梯度消失问题

③ ReLU: 优点: 不会出现饱和

计算代价小

缺点: 非0中心分布

dead ReLU

④ Leaky ReLU: $\max(0.1x, x)$ 变种: PReLU = $\max(px, x)$

优点: 不会出现饱和

计算代价小

收敛速度快

不会dead

讲师经验之谈!

TLDR: In practice:

- Use ReLU. Be careful with your learning rates
- Try out Leaky ReLU / Maxout / ELU
- Try out tanh but don't expect much
- Don't use sigmoid

1.2.2 数据预处理

Step 1: Preprocess the data

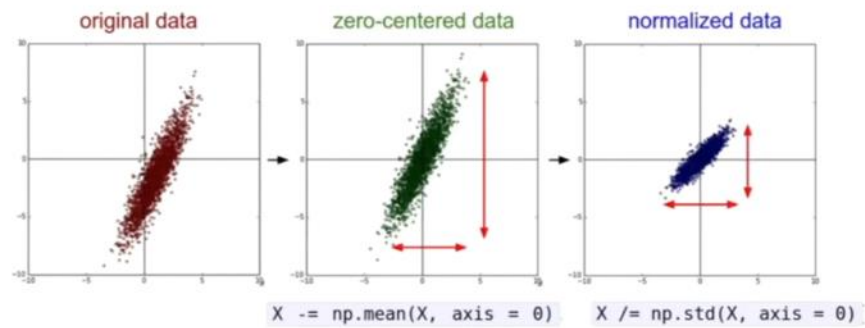


图6.1.3 数据预处理

1. 数据分布中心化处理
2. 归一化处理
3. PCA 主成分分析
降维, 分析数据相关性
4. 白化

2. 批量归一化

《Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift》

背景：在网络训练的过程中，参数的更新会导致网络的各层输入数据的分布不断变化，那么各层在训练的过程中就需要不断的改变以适应这种新的数据分布，从而造成网络训练困难，收敛变慢（而且网络越深越难），在论文中这个问题被称为“Internal Covariate Shift”。为了解决这个问题出现了批量归一化的算法，他对每一层的输入进行归一化，保证每层的输入数据分布是稳定的，从而加速训练。

Batch Normalization

[Ioffe and Szegedy, 2015]

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1, \dots, x_m\}$;
Parameters to be learned: γ, β

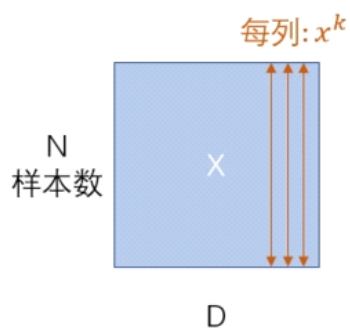
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

- Improves gradient flow through the network
- Allows higher learning rates
- Reduces the strong dependence on initialization
- Acts as a form of regularization in a funny way, and slightly reduces the need for dropout, maybe

Stanford

$$\hat{x}^{(k)} = \frac{x^k - E[x^k]}{\sqrt{Var[x^k]}}$$



计算数据每个维度（每列）的均值和方差：

$$\text{第}k\text{维的平均值 } E[x^k] = \frac{1}{N} \sum_{j=1}^N x_j$$

$$\text{第}k\text{维的标准差 } \sqrt{Var[x^k]} = \sqrt{\frac{1}{N} \sum_{j=1}^N (x_j - E[x^k])^2}$$

再用上面的公式对每一维进行归一化

BN: <https://blog.csdn.net/hjimce/article/details/50866313>