第三讲 全连接网络 Lecture 3 Fully-connected Neural Networks

明玉瑞 Yurui Ming yrming@gmail.com

声明 Disclaimer

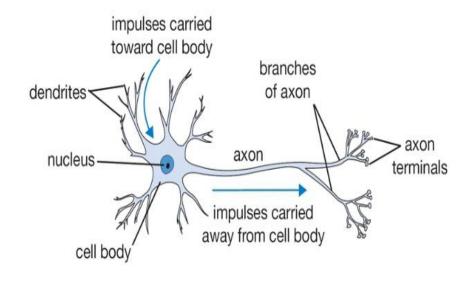
本讲义在准备过程中由于时间所限,所用材料来源并未规范标示引用来源。所引材料仅用于教学所用,作者无意侵犯原著者之知识版权,所引材料之知识版权均归原著者所有;若原著者介意之,请联系作者更正及删除。

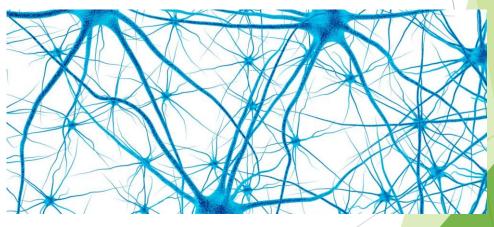
The time limit during the preparation of these slides incurs the situation that not all the sources of the used materials (texts or images) are properly referenced or clearly manifested. However, all materials in these slides are solely for teaching and the author is with no intention to infringe the copyright bestowed on the original authors or manufacturers. All credits go to corresponding IP holders. Please address the author for any concern for remedy including deletion.

脑科学启发 Brain inspiration

▶ 全连接网络基于神经元处理信息的方式及其组网方式

Fully-connected neural networks are inspired by the way neurons work and their connections into network





数学模型

Mathematical Modelling

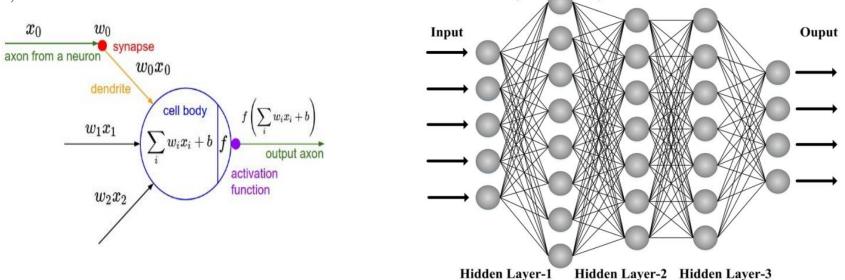
> 对神经元及其构成网络的数学建模

Mathematical modelling of the neurons and corresponding networks

▶ 为建模方便,引入了额外的术语,如输入层,隐含层,输出层等

Terminologies are introduced for modelling, such as input layer, hidden layer, output

layer, etc.



 W_{lm}

记号

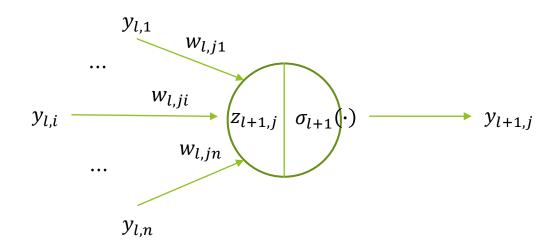
Denotation

- ▶ 第l层的第i个神经元输出记为 $y_{l,i}$ The output of i-th neuron at level l is denoted by $y_{l,i}$
- ▶ 第l层的第i个神经元到第l+1层的第j个神经元记为 $w_{l,ji}$ The weight from the i-th neuron at level l to the j-th neuron at level l+1 is denoted by $w_{l,ji}$
- b 由第l层的神经元到第l+1层的第j个神经元的加权和记为 $z_{l+1,j}$ The intermediate result for the j-th neuron at level l+1 is denoted by $z_{l+1,j}$

记号

Denotation

▶ 用于建模的数学式子展示 Illustration of mathematical modelling formulas



$$z_{l+1,j} = \sum_{i=1}^{n} y_{l,i} \cdot w_{l,ji}$$

$$y_{l+1,j} = \sigma_{l+1}(z_{l+1,j}) = \frac{1}{1 + e^{-z_{l+1,j}}}$$
(sigmoid function)

feed-forward propagation

▶ 可以逐节点算

Node-based computation

$$y_{6} = x_{6} y_{7} = x_{7}$$

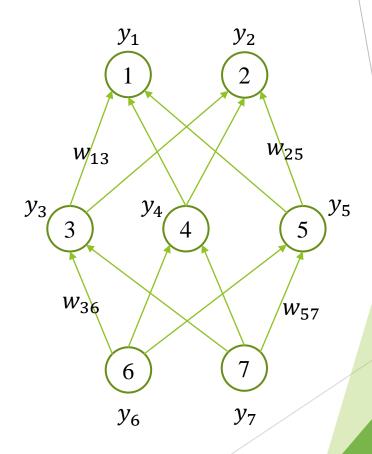
$$y_{3} = \sigma(w_{36} \cdot y_{6} + w_{37} \cdot y_{7})$$

$$y_{4} = \sigma(w_{46} \cdot y_{6} + w_{47} \cdot y_{7})$$

$$y_{5} = \sigma(w_{56} \cdot y_{6} + w_{57} \cdot y_{7})$$

$$y_{1} = \sigma(w_{13} \cdot y_{3} + w_{14} \cdot y_{4} + w_{15} \cdot y_{5})$$

$$y_{2} = \sigma(w_{23} \cdot y_{3} + w_{24} \cdot y_{4} + w_{25} \cdot y_{5})$$



feed-forward propagation

▶ 实践中会是矩阵的形式 Matrix form in practice

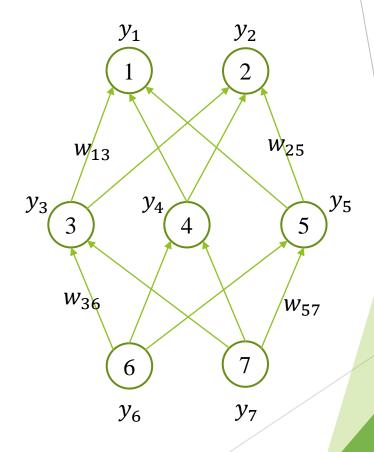
$$y_{3} = \sigma(w_{36} \cdot y_{6} + w_{37} \cdot y_{7})$$

$$y_{4} = \sigma(w_{46} \cdot y_{6} + w_{47} \cdot y_{7})$$

$$y_{5} = \sigma(w_{56} \cdot y_{6} + w_{57} \cdot y_{7})$$

$$\begin{bmatrix} y_{3} \\ y_{4} \\ y_{5} \end{bmatrix} = \sigma(\begin{bmatrix} w_{36} & w_{37} \\ w_{37} & w_{37} \\ w_{37} & w_{37} \end{bmatrix} \begin{bmatrix} y_{6} \\ y_{7} \end{bmatrix})$$

$$\vec{y}_{l+1} = \sigma(W\vec{y}_l)$$

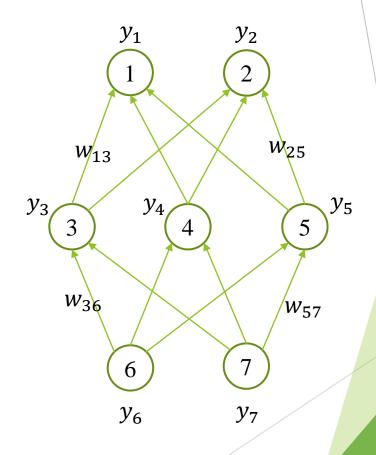


feed-forward propagation

▶ 实践中会是矩阵的形式 Matrix form in practice

$$y_1 = \sigma(w_{13} \cdot y_3 + w_{14} \cdot y_4 + w_{15} \cdot y_5)$$
$$y_2 = \sigma(w_{23} \cdot y_3 + w_{24} \cdot y_4 + w_{25} \cdot y_5)$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \sigma \left(\begin{bmatrix} w_{13} & w_{14} & w_{15} \\ w_{23} & w_{24} & w_{25} \end{bmatrix} \begin{bmatrix} y_3 \\ y_4 \\ y_5 \end{bmatrix} \right)$$
$$\vec{y}_{l+1} = \sigma(W\vec{y}_l)$$



feed-forward propagation

▶ 进一步考虑阀限值

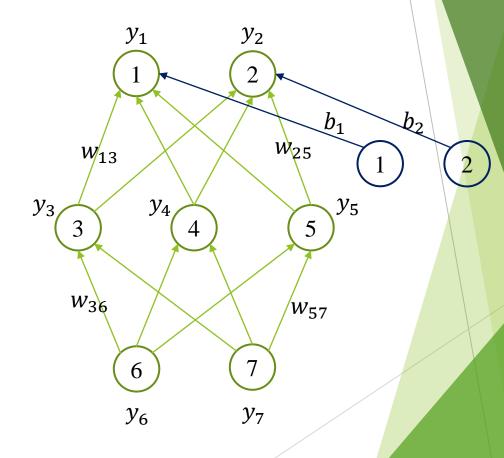
Put bias into consideration

$$y_{1} = \sigma(w_{13} \cdot y_{3} + w_{14} \cdot y_{4} + w_{15} \cdot y_{5} + b_{1})$$

$$y_{2} = \sigma(w_{23} \cdot y_{3} + w_{24} \cdot y_{4} + w_{25} \cdot y_{5} + b_{2})$$

$$\begin{bmatrix} y_{1} \\ y_{2} \end{bmatrix} = \sigma \begin{pmatrix} \begin{bmatrix} w_{13} & w_{14} & w_{15} \\ w_{23} & w_{24} & w_{25} \end{bmatrix} \begin{bmatrix} y_{3} \\ y_{4} \\ y_{5} \end{bmatrix} + \begin{bmatrix} b_{1} \\ b_{2} \end{bmatrix} \end{pmatrix}$$

$$\vec{y}_{l+1} = \sigma(W\vec{y}_{l} + \vec{b})$$



feed-forward propagation

▶ 进一步考虑批量的情况

Put batch into consideration

$$y_{1}^{1} = \sigma(w_{13} \cdot y_{3}^{1} + w_{14} \cdot y_{4}^{1} + w_{15} \cdot y_{5}^{1} + b_{1})$$

$$y_{2}^{1} = \sigma(w_{23} \cdot y_{3}^{1} + w_{24} \cdot y_{4}^{1} + w_{25} \cdot y_{5}^{1} + b_{2})$$

$$y_{2}^{2} = \sigma(w_{23} \cdot y_{3}^{2} + w_{24} \cdot y_{4}^{2} + w_{25} \cdot y_{5}^{2} + b_{2})$$
...
$$y_{1}^{N} = \sigma(w_{13} \cdot y_{3}^{N} + w_{14} \cdot y_{4}^{N} + w_{15} \cdot y_{5}^{N} + b_{1})$$

$$y_{2}^{2} = \sigma(w_{23} \cdot y_{3}^{2} + w_{24} \cdot y_{4}^{2} + w_{25} \cdot y_{5}^{2} + b_{2})$$
...
$$y_{1}^{N} = \sigma(w_{13} \cdot y_{3}^{N} + w_{14} \cdot y_{4}^{N} + w_{15} \cdot y_{5}^{N} + b_{1})$$

$$y_{2}^{N} = \sigma(w_{23} \cdot y_{3}^{N} + w_{24} \cdot y_{4}^{N} + w_{25} \cdot y_{5}^{N} + b_{2})$$

$$\begin{bmatrix} y_1^1 & y_2^2 & y_2^3 & y_2^4 \\ y_2^1 & y_2^2 & y_2^3 & y_2^4 \end{bmatrix} = \sigma \left(\begin{bmatrix} w_{13} & w_{14} & w_{15} \\ w_{23} & w_{24} & w_{25} \end{bmatrix} \begin{bmatrix} y_3^1 & y_3^2 & y_3^3 & y_3^4 \\ y_4^1 & y_4^2 & y_4^3 & y_4^4 \\ y_5^1 & y_5^2 & y_5^3 & y_5^4 \end{bmatrix} + \begin{bmatrix} b_1 & b_1 & b_1 & b_1 \\ b_2 & b_2 & b_2 & b_2 \end{bmatrix} \right)$$

$$Y_{l+1} = \sigma(WY_l + B)$$

feed-forward propagation

▶ 工程习惯中,批量大小用N表示,作为第一个维度;特征列大小用C表示,作为第二个维度,因此需要转置一下

From the engineering convention perspective, the batch size is denoted by N as the first dimension, while the size of feature column is denoted by C as the second dimension, so the formula deducted in the last slide needs to be transposed.

$$\begin{bmatrix} y_1^1 & y_2^2 & y_2^3 & y_2^4 \\ y_2^1 & y_2^2 & y_2^3 & y_2^4 \end{bmatrix} = \sigma \left(\begin{bmatrix} w_{13} & w_{14} & w_{15} \\ w_{23} & w_{24} & w_{25} \end{bmatrix} \begin{bmatrix} y_3^1 & y_3^2 & y_3^3 & y_3^4 \\ y_4^1 & y_4^2 & y_4^3 & y_4^4 \\ y_5^1 & y_5^2 & y_5^3 & y_5^4 \end{bmatrix} + \begin{bmatrix} b_1 & b_1 & b_1 & b_1 \\ b_2 & b_2 & b_2 & b_2 \end{bmatrix} \right)$$

$$\begin{bmatrix} y_1^1 & y_2^1 \\ y_1^2 & y_2^2 \\ y_1^3 & y_2^3 \\ y_1^4 & y_2^4 \end{bmatrix} = \sigma \begin{pmatrix} \begin{bmatrix} y_3^1 & y_4^1 & y_5^1 \\ y_3^2 & y_4^2 & y_5^2 \\ y_3^3 & y_4^3 & y_5^3 \\ y_3^4 & y_4^4 & y_5^4 \end{bmatrix} \begin{bmatrix} w_{13} & w_{23} \\ w_{14} & w_{24} \\ w_{15} & w_{25} \end{bmatrix} + \begin{bmatrix} b_1 & b_2 \\ b_1 & b_2 \\ b_1 & b_2 \\ b_1 & b_2 \end{bmatrix}$$

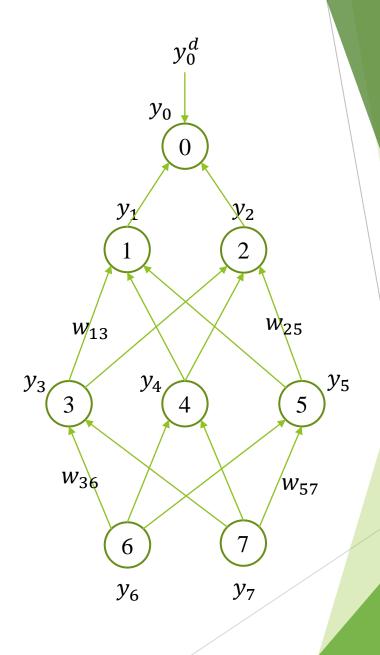
$$Y_{l+1} = \sigma(Y_lW + B)$$

feed-forward propagation

▶ 计算损失,这里以回归问题为例
Loss computation, exemplified by regression problem

$$E = \frac{1}{2} \left(y_0 - y_0^d \right)^2$$

$$E = \frac{1}{2N} \sum_{i=1}^{N} (y_0^i - y_0^d)^2$$

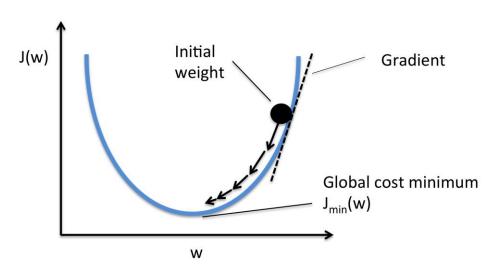


反向传播算法 Back-propagation algorithm

- ▶ 反向传播算法又称错误传播或delta学习法则
 An alternative for back propagation is called error propagation or delta learning rule
- ▶ 为方便推导,做下面一些假设
 The following assumptions are made to ease the deduction
 - $ightharpoonup 若层不存在歧义,则<math>w_{l,ji}$ 简记为 w_{ji} If no ambiguity from the layer perspective, $w_{l,ji}$ is denoted by w_{ji}
 - ▶ 对于输出层的第i个神经元,其对应的真实标签记为 y_i^d For output from the i-th neurons at the final layer L, the ground-truth is denoted by $y_{L,i}^d$, or y_i^d

Back-propagation algorithm

- ▶ 反向传播算法基本思想基于梯度下降
 Back-propagation algorithm is based on gradient descent algorithm
- ▶ 反向传播算法的实施依赖于求导的链式法则
 Implementation of back-propagation algorithm relies on the chain rules of derivatives



$$E = \frac{1}{2} (y_i - y_i^d)^2$$

$$w_{ij}(t+1) = w_{ij}(t) + \eta \left(-\frac{\partial E}{\partial w_{ij}} \right)$$

$$\frac{\partial E}{\partial w_{ij}} = (y_i - y_i^d) \frac{\partial y_i}{\partial w_{ij}} = (y_i - y_i^d) \frac{\partial y_i}{\partial z_i} \frac{\partial z_i}{\partial w_{ij}} = (y_i - y_i^d) y_i' y_j$$

$$w_{ij}(t+1) = w_{ij}(t) - \eta (y_i - y_i^d) y_i' y_j$$

Back-propagation algorithm

▶ 简单起见, 我们顺序标识神经元节点进行推导

For simplicity, neurons are labelled sequentially

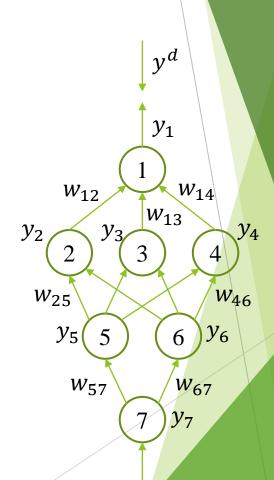
$$E = \frac{1}{2} (y_1 - y^d)^2$$
 $\Rightarrow \frac{\partial E}{\partial y_1} = (y_1 - y^d)$

$$\frac{\partial E}{\partial w_{12}} = \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial w_{12}} = (y_1 - y^d) \frac{\partial y_1}{\partial z_1} \frac{\partial z_1}{\partial w_{12}} = (y_i - y_i^d) y_1' y_2$$

$$\delta_0 \triangleq (y_1 - y^d)$$
 $\Rightarrow \frac{\partial E}{\partial w_{12}} = \delta_0 y_1' y_2$

$$\delta_1 \triangleq \delta_0 y_1' \qquad \Rightarrow w_{12}(t+1) = w_{12}(t) - \eta \delta_1 y_2$$
$$w_{13}(t+1) = w_{13}(t) - \eta \delta_1 y_3$$

$$w_{14}(t+1) = w_{14}(t) - \eta \delta_1 y_4$$



Back-propagation algorithm

$$\frac{\partial E}{\partial w_{25}} = \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial w_{25}} = (y_1 - y^d) \frac{\partial y_1}{\partial z_1} \frac{\partial z_1}{\partial w_{25}} = (y_i - y_i^d) y_1' \times$$

$$\frac{\partial}{\partial w_{25}} [w_{12} y_2 + w_{13} y_3 + w_{14} y_4] = \delta_1 w_{12} \frac{\partial y_2}{\partial w_{25}}$$

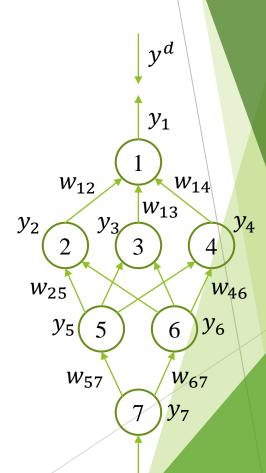
$$= \delta_1 w_{12} \frac{\partial y_2}{\partial z_2} \frac{\partial z_2}{\partial w_{25}} = \delta_1 w_{12} y_2' \frac{\partial z_2}{\partial w_{25}}$$

$$= \delta_1 w_{12} \frac{\partial y_2}{\partial z_2} \frac{\partial}{\partial w_{25}} [w_{25} y_5 + w_{26} y_6] = \delta_1 w_{12} y_2' y_5$$

$$\delta_2 \triangleq \delta_1 w_{12} y_2' \qquad \Rightarrow \frac{\partial E}{\partial w_{25}} = \delta_2 y_5$$

$$w_{25} (t+1) = w_{25} (t) - \eta \delta_2 y_5$$

$$w_{26} (t+1) = w_{26} (t) - \eta \delta_2 y_6$$



Back-propagation algorithm

Similarly, we have

$$\delta_{3} \triangleq \delta_{1}w_{13}y_{3}'$$

$$\delta_{4} \triangleq \delta_{1}w_{14}y_{4}'$$

$$\delta_{5} \triangleq (\delta_{2}w_{25} + \delta_{3}w_{35} + \delta_{4}w_{45})y_{5}'$$

$$\delta_{6} \triangleq (\delta_{2}w_{26} + \delta_{3}w_{36} + \delta_{4}w_{46})y_{6}'$$

$$w_{35}(t+1) = w_{35}(t) - \eta \delta_{3}y_{5}$$

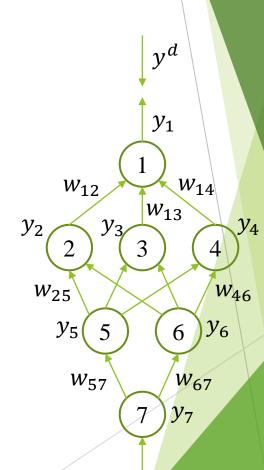
$$w_{36}(t+1) = w_{36}(t) - \eta \delta_{3}y_{6}$$

$$w_{45}(t+1) = w_{45}(t) - \eta \delta_{4}y_{5}$$

$$w_{46}(t+1) = w_{46}(t) - \eta \delta_{4}y_{6}$$

$$w_{57}(t+1) = w_{57}(t) - \eta \delta_{5}y_{7}$$

$$w_{67}(t+1) = w_{67}(t) - \eta \delta_{6}y_{7}$$



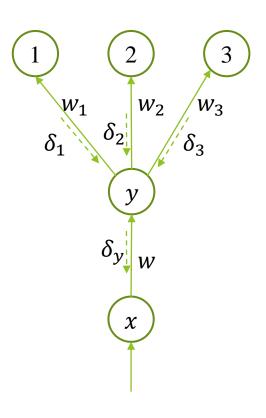
Back-propagation algorithm

一般性规则

General rule

$$\delta_y \triangleq (\delta_1 w_1 + \delta_2 w_2 + \delta_3 w_3) y'$$

$$w(t+1) = w(t) - \eta \delta_y x$$



实践过程 Practice

▶ 交替执行以下过程

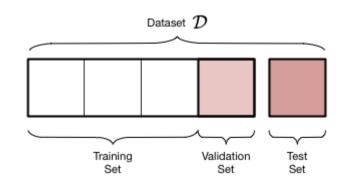
Perform the following procedures in an iterative way:

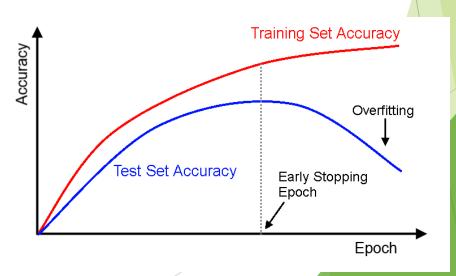
▶ 调优网络配置,训练网络

Decide the neural network configuration and train the neural network

▶ 进行(交叉)验证,评估网络设计及性能

(Cross-)validate to evaluate the design and performance





习题 Problems

阅读相关文献,解释如何根据在训练过程中,根据模型在训练集与验证集上的性 能来判断欠拟合与过拟合。

Explain how to interpret the underfitting and overfitting phenomena during the training process by inspecting the loss and/or accuracy on training set and test set, by doing literature survey.

2. 请阅读文献, 查阅异或 (XOR) 问题的历史, 了解其对神经网络发展的影响, 并 设计一个全连接网络解决这个问题。

By doing literature review to understand the XOR problems in the history of neural network evolvement, and its impact on the declining of neural networks. Design a multilayer fully-connected neural network to solve the XOR problem.

习题

Problems

- 3. 在Slide 4的公式中, σ称为激活函数 σ is called activation function in the slide 4
 - 请举例说明你所知道的激活函数; Enumerate some examples of activation functions.
 - 请证明若激活函数为恒等变换(Identity), 即 $\sigma(x) = x$, 则多层神经网络等价于 线性变换;
 - Proof that if identity function is adopted as activation function, then multi-perceptron network is equivalent to linear transformation.
 - 对于一些特定的激活函数,例如Sigmoid函数,在实际应用中可能会有梯度消失 (Gradient vanishing)问题,请解释,并证明其最大梯度值为0.25;
 - For some activation functions, such as sigmoid function, there could exist gradient vanishing problem in practice. Please explain it and prove that the maximal gradient is 0.25

习题 Problems

• 单调性是通常认为是激活函数的一个必要条件,请阅读最新文献,说明这个条件 非必要, 即非单调函数有可能在某些情况下表现更好; 并尝试从神经科学的角度 解释这件事。

Monotone is attributed to the necessity for the candidature of activation functions. By doing literature survey to illustrate that non-monotonic function could have better performance, which dis-prove the proclamation. Find explanations from the neuroscientific perspective.

4. 偏差和方差是机器学习模型的两个重要属性。偏差是为简化模型优化过程而做出的假 设,方差是目标函数可以输出的值的范围的度量。请说明如果偏差与方差过高,都有 哪些危害。

Bias and variance are two factors measuring the quality of ML models. Bias is introduced due to assumption for model simplification; while variance measures ranges of model outputs. By doing literature review to clarify the harmfulness of high bias and high variance.