

1. Read **Deep Learning: An Introduction for Applied Mathematicians**. Consider a network as defined in (3.1) and (3.2). Assume that $n_L = 1$, find an algorithm to calculate $\nabla a^{[L]}(x)$.

1

Goal: $\nabla_x H(x) = \left[\frac{\partial H}{\partial x_1}, \frac{\partial H}{\partial x_2}, \dots, \frac{\partial H}{\partial x_n} \right]^T$

Input: $a^{[0]} = \vec{x} \in \mathbb{R}^n$, Having L layer

where $z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}$, $a^{(l)} = \sigma(z^{(l)})$

→ Outcome is $H(x) = a^{[L]} \in \mathbb{R}^1$

<sol>

1° last layer, $a^{[L]} = \sigma(z^{[L]}) = H(x) \Rightarrow \delta^{[L]} = \frac{\partial H}{\partial z^{[L]}} = \sigma'(z^{[L]})$

For each layer, we have $\delta^{(l)} = \frac{\partial H}{\partial z^{(l)}} = \frac{\partial z^{(l+1)T}}{\partial z^{(l)}} \frac{\partial H}{\partial z^{(l+1)}} \odot \sigma'(z^{(l)})$

Moreover, we have $z^{(l+1)} = W^{(l+1)} a^{(l)} + b^{(l+1)}$, $a^{(l)} = \sigma(z^{(l)})$

① $\frac{\partial z^{(l+1)}}{\partial a^{(l)}} = W^{(l+1)}$ ② $\frac{\partial a^{(l)}}{\partial z^{(l)}} = \text{diag}(\sigma'(z^{(l)}))$

→ $\delta^{(l)} = (W^{(l+1)})^T \delta^{(l+1)} \odot \sigma'(z^{(l)})$

2° By Chain Rule: $\frac{\partial H}{\partial a^{[0]}} = \frac{\partial z^{[1]T}}{\partial a^{[0]}} \frac{\partial H}{\partial z^{[1]}} = (W^{[1]})^T \delta^{[1]}$

Therefore, we have, $\nabla_x H(x) = \frac{\partial H}{\partial x} = (W^{[1]})^T \delta^{[1]}$

→ $\nabla_x a(x) = (W^{[1]})^T \delta^{[1]}$ ■

2. There are unanswered questions during the lecture, and there are likely more questions we haven't covered. Take a moment to think about them and write them down here.

☆在做HW2的時候我遇到了一些問題,首先是 activation function 真的能逼近一個 continuous function 嗎? 令一個問題是我該如果驗證我的 model 與我的原函數很近似,所以下面有 ① Universal Approximation Theorem 的討論。
② Using Training sets & Validation sets

▷ Universal Approximation Theorem (通用逼近定理)

又叫万能近似定理。

→ 一个单隐藏层的 feedforward Neural Network, 只要 hidden layer 神经元 enough, 且激活的函数 content 一些条件 (sigmoid, ReLU, tanh), 就能逼近任意的连续函数, (在一个有界区间, 任意的误差)。

Assume $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuous function.

$$F(x) = \sum_{i=1}^m \alpha_i \sigma(w_i x + b_i) \leftarrow \text{存在一个单hidden layer 的 neural network.}$$

$$\text{s.t. } \forall \epsilon > 0, |f(x) - F(x)| < \epsilon, \forall x \in \mathbb{R}^n \text{ 定义域}$$

▷ Early Stopping (避免 ML 过拟合)

(用来更新权重的训练集 model 在看过资料的表现)

▷ 一开始有一个完整的数据: Input 和 labels $D = \{(x_i, y_i)\}_{i=1}^N$

① 训练集损失函数:

$$\mathcal{L}_{\text{train}} = \frac{1}{|D_{\text{train}}|} \sum_{(x,y) \in D_{\text{train}}} \ell(f_\theta(x), y)$$

② 验证集损失函数:

$$\mathcal{L}_{\text{val}} = \frac{1}{|D_{\text{val}}|} \sum_{(x,y) \in D_{\text{val}}} \ell(f_\theta(x), y)$$

Train ML model =

① (Training set, D_{train})

② (Validation set, D_{val})

③ (Test set, D_{test}) (可选)

假设有 1000 条 Data, 可能切成:

① D_{train} : 700 条 (70%)
② D_{val} : 200 条 (20%)
③ D_{test} : 100 条 (10%)

这是与参数更新的。

没有参数训练的验证集

→ In Early stopping: we calculate every epoch: model 在看过资料的表现。

(1) $\mathcal{L}_{\text{train}}(\theta^{(t)})$: 训练的误差 (通常一直下降)

(2) $\mathcal{L}_{\text{val}}(\theta^{(t)})$: 验证误差, 下降到某处后会上升。

⊗ Early stopping 就观察 \mathcal{L}_{val} , 并在恶化前停下来!

→ 在第 t 个 epoch 训练后, 得到参数 θ^t

→ Early stopping 的目标就是 find: $t^* = \arg \min_{t \leq T} \mathcal{L}_{\text{val}}(\theta^{(t)})$ (最大训练回合数)

→ Patience (耐心值) (P)

$\mathcal{L}_{\text{val}}(\theta^{(t)}) > \min_{s \leq t} \mathcal{L}_{\text{val}}(\theta^{(s)})$, 连续 P 次, 则停止 Train.

Lately, 保存最佳 Model: $\theta^* = \theta^{(t^*)}$, $t^* = \arg \min_{t \leq T} \mathcal{L}_{\text{val}}(\theta^{(t)})$.