

Deep Learning!

Meichen Lu (ml574@cam.ac.uk)

April 19, 2018

Contents

1	Classic neural network	1
1.1	Components	1
1.1.1	Structure	1
1.1.2	Neurons	1
1.1.3	Log likelihood	2
1.1.4	Notations	2
1.2	Vectorisation	2
1.3	Back-propagation	3
1.3.1	Weights initialization	3
1.3.2	Loss function	3
1.3.3	Gradient descend	3

1 Classic neural network

1.1 Components

Intuition: 'automatic' feature engineering and discovery! Each unit is like a neuron in the brain.

1.1.1 Structure

- Input features: input layer
- Output neuron: output layer
- Everything in the middle: hidden layers (with hidden units)

1.1.2 Neurons

Each neuron 'evaluates' the input according to a certain features and is represented as a function. It is called **activation function**. We can break down the function into two distinct computations: (1) $z = w^T x + b$ and (2) $a = g(z)$. Examples of $g(z)$ are

- Sigmoid: $g(z) = \frac{1}{1+e^{-z}}$
- ReLu: $g(z) = \max(z, 0)$
- tanh: $g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

1.1.3 Log likelihood

The probability from sigmoid function in the final layer is the log-likelihood function

$$\sum_{i=1}^m \left(y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right) \quad (1)$$

1.1.4 Notations

- input: $\mathbf{x} = [x_1, x_2, \dots]$
 - Sample is noted as $\mathbf{x}^{(i)}$
- Linear combination of inputs: \mathbf{z}
- Neuron output: \mathbf{a}
- Layer notation: $\mathbf{z}^{[1]}$ for first layer
- Final output: y
- Predicted output: \hat{y}

Putting it all together, the first first hidden unit in the first hidden layer performs the following computations:

$$z_1^{[1]} = \mathbf{W}_1^{[1]T} \mathbf{x} + b_1^{[1]} \quad \text{and} \quad a_1^{[1]} = g(z_1^{[1]}) \quad (2)$$

1.2 Vectorisation

- Fully training set input:

$$X = \begin{pmatrix} \begin{array}{c} | \\ \mathbf{x}^{(1)} \\ | \end{array} & \begin{array}{c} | \\ \mathbf{x}^{(2)} \\ | \end{array} & \dots & \begin{array}{c} | \\ \mathbf{x}^{(m)} \\ | \end{array} \end{pmatrix}$$

- Fully stacked weight matrix

$$W^{[1]} = \begin{bmatrix} - & \mathbf{W}_1^{[1]T} & - \\ - & \mathbf{W}_2^{[1]T} & - \\ & \vdots & \\ - & \mathbf{W}_{n_1}^{[1]T} & - \end{bmatrix}$$

- Fully projection in the first layer from the inputs X is thus

$$Z^{[1]} = \begin{pmatrix} \begin{array}{c} | \\ \mathbf{z}^{1} \\ | \end{array} & \begin{array}{c} | \\ \mathbf{z}^{[1](2)} \\ | \end{array} & \dots & \begin{array}{c} | \\ \mathbf{z}^{[1](m)} \\ | \end{array} \end{pmatrix} = W^{[1]}X + b^{[1]}$$

1.3 Back-propagation

1.3.1 Weights initialization

Usually, we randomly initialize the parameters to small values (e.g., normally distributed around zero; $\mathcal{N}(0, 0.1)$) In practice, it turns out there is something better than random initialization. It is called Xavier/He initialization and initializes the weights:

$$w^{[\ell]} \sim \mathcal{N}\left(0, \sqrt{\frac{2}{n^{[\ell]} + n^{[\ell-1]}}}\right)$$

1.3.2 Loss function

With regularisation

$$J(W, b) = \left[\frac{1}{m} \sum_{i=1}^m J(W, b; \mathbf{x}^{(i)}, y^{(i)}) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} \left(W_{ji}^{[l]} \right)^2 \quad (3)$$

The weight decay parameter λ controls the relative importance of the two terms.

1.3.3 Gradient descend

Update rules for gradient descent

$$W_{ij}^{[l]} = W_{ij}^{[l]} - \alpha \frac{\partial}{\partial W_{ij}^{[l]}} J(W, b) \quad (4)$$

$$b_i^{[l]} = b_i^{[l]} - \alpha \frac{\partial}{\partial b_i^{[l]}} J(W, b) \quad (5)$$

Or with regularization

$$\frac{\partial}{\partial W_{ij}^{[l]}} J(W, b) = \left[\frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial W_{ij}^{[l]}} J(W, b; \mathbf{x}^{(i)}, y^{(i)}) \right] + \lambda W_{ij}^{[l]} \quad (6)$$

$$\frac{\partial}{\partial b_i^{[l]}} J(W, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial b_i^{[l]}} J(W, b; \mathbf{x}^{(i)}, y^{(i)}) \quad (7)$$

Stochastic gradient descend Calculating the influence on gradient from a single example, using the cost function:

$$J(W, b; \mathbf{x}, y) = \frac{1}{2} \|h_{W,b}(\mathbf{x}) - y\|^2 \quad (8)$$

1. Perform a feedforward pass up to the output layer L_{n_l}
2. For each output unit i in the output layer n_l , compute

$$\delta_i^{[n_l]} = \frac{\partial}{\partial z_i^{[n_l]}} \frac{1}{2} \|y - h_{W,b}(\mathbf{x})\|^2 = -(y_i - a_i^{[n_l]}) \cdot f'(z_i^{[n_l]}) \quad (9)$$

3. For $l = n_l - 1, n_l - 2, \dots, 2$
 For each node i in layer l , set

$$\delta_i^{[l]} = \left(\sum_{j=1}^{s_{l+1}} W_{ji}^{[l]} \delta_j^{[l+1]} \right) f'(z_i^{[l]}) \quad (10)$$

4. Compute the partial derivatives:

$$\frac{\partial}{\partial W_{ij}^{[l]}} J(W, b; \mathbf{x}, y) = a_j^{[l]} \delta_i^{[l+1]} \quad (11)$$

$$\frac{\partial}{\partial b_i^{[l]}} J(W, b; \mathbf{x}, y) = \delta_i^{[l+1]}. \quad (12)$$

Or in the vectorised form

1. Feedforward
2. $\delta^{[n_l]} = -(y - a^{[n_l]}) \circ f'(z^{[n_l]})$
3. $\delta^{[l]} = ((W^{[l]})^T \delta^{[l+1]}) \circ f'(z^{[l]})$
4. $\nabla_{W^{[l]}} J(W, b; \mathbf{x}, y) = \delta^{[l+1]} (a^{[l]})^T$

\circ is the Hadamard product, i.e., element-wise multiplication