



Neural Network

Li Ju

22/06 - 03/07

Content

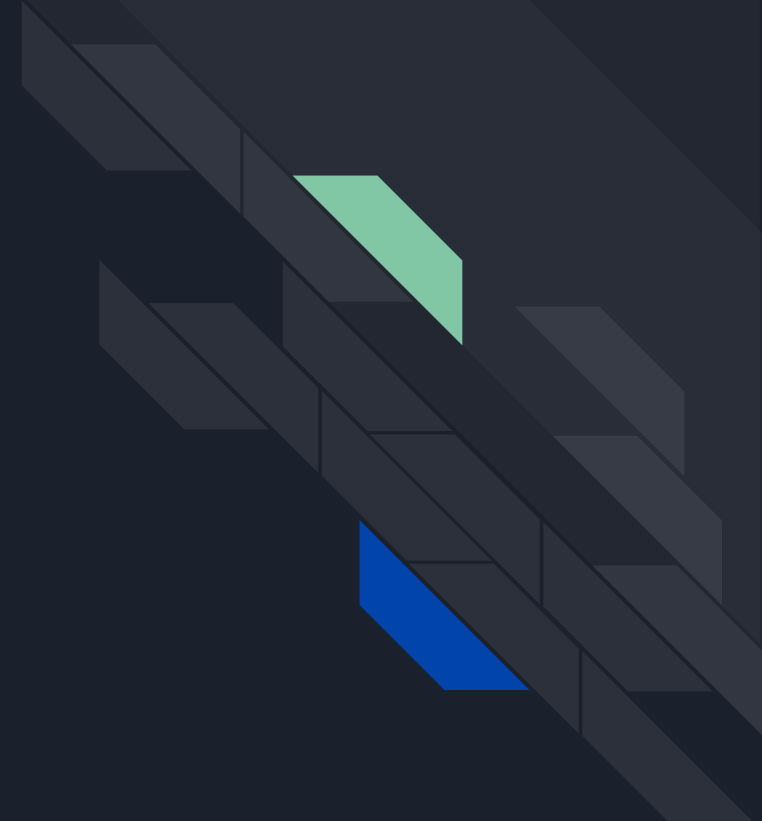
Overview

Keras implementation

C implementation

Activation functions/Loss functions ...

Chosen datasets for test

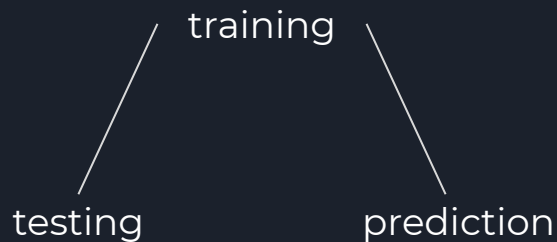




Overview

Classification/regression model:

Input --(Neural Network)--> Output



Advantages:

1. Flexibility
2. Less assumption
3. Generalizability

Disadvantages:

1. Costly computation
2. Lack of interpretability
3. Structure flexibility



Keras implementation

training_x matrix (features)



input_dim = #features, output_dim = #nodes, AF = "relu"



input_dim = #upper_nodes, output_dim = #nodes, AF = "relu"



input_dim = #upper_nodes, output_dim = #categories, AF = "softmax"



output_y matrix ~ target_y: Loss function -> Error

Gradient descending:
Minimize error

Cross validation:
Avoid overfitting

Epoch:
#training times through
entire dataset

Batch size:
#samples for each iteration



However...

Training mechanism?

We need mathematical details!

A detailed close-up of a printed circuit board (PCB) is shown in the bottom-left corner, partially obscured by a dark blue diagonal overlay. The PCB features various electronic components, including a large black electrolytic capacitor labeled '132400', several resistors with values like '475', and other smaller components labeled 'R4', 'R5', 'R3', and 'C4'. A soldering iron tip is visible, working on a component on the board.

C implementation for a single-hidden-layer network

Scalar-based

Matrix-based

Notation:

$X_{m \times n}$: Input matrix with m samples and n features

$Y_{m \times t}$: Target matrix with m samples and t predicted output

Q^T : Transpose of matrix Q

L : Output layer

n^a : number of nodes in layer a

$W_{n^{a-1} \times n^a}^a$: Weight matrix between upper layer $a - 1$ and lower layer a , the dimension of which is $n^{a-1} \times n^a$

$B_{1 \times n^a}^a$: Bias matrix(vector) between upper layer $a - 1$ and lower layer a , the dimension of which is $1 \times n^a$

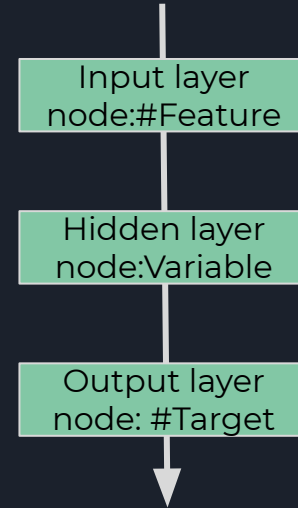
$Z_{m \times n^a}^a$: Weighted input for layer a , the dimension of which is $m \times n^a$

$A_{m \times n^a}^a$: Output matrix of layer $a - 1$, or, input matrix of layer a

$\sigma^a(Z)$: Activation function of layer a mapping matrix Z element-wise

E : Error calculated by loss function

$\frac{\partial Y}{\partial X}$: Differential of variable Y on variable X



Initialization:

Let all elements of every **W** and **B** matrices be a random number within -1 and 1;

Forward propagation:

$$\mathbf{A}_{m \times t}^{L+1} = \sigma^L((\dots \sigma^2(\sigma^1(\mathbf{X}_{m \times n} \cdot \mathbf{W}_{n \times n^1}^1 + \mathbf{B}_{1 \times n^1}^1) \mathbf{W}_{n^1 \times n^2}^2 + \mathbf{B}_{1 \times n^2}^2) \dots) \dots \mathbf{W}_{n^{L-1} \times t}^L + \mathbf{B}_{1 \times t}^L)$$

$\mathbf{X}_{m \times n}$: Input matrix with m samples and n features

$\mathbf{Y}_{m \times t}$: Target matrix with m samples and t predicted output

\mathbf{Q}^T : Transpose of matrix \mathbf{Q}

L : Output layer

n^a : number of nodes in layer a

$\mathbf{W}_{n^{a-1} \times n^a}^a$: Weight matrix between upper layer $a - 1$ and lower layer a , the dimension of which is $n^{a-1} \times n^a$

$\mathbf{B}_{1 \times n^a}^a$: Bias matrix(vector) between upper layer $a - 1$ and lower layer a , the dimension of which is $1 \times n^a$

$\mathbf{Z}_{m \times n^a}^a$: Weighted input for layer a , the dimension of which is $m \times n^a$

$\mathbf{A}_{m \times n^a}^a$: Output matrix of layer $a - 1$, or, input matrix of layer a

$\sigma^a(\mathbf{Z})$: Activation function of layer a mapping matrix \mathbf{Z} element-wise

E : Error calculated by loss function

$\frac{\partial Y}{\partial X}$: Differential of variable Y on variable X

Error:

$$E = \frac{\sum (\mathbf{A}_{m \times t}^{L+1} - \mathbf{Y}_{m \times t}) \odot (\mathbf{A}_{m \times t}^{L+1} - \mathbf{Y}_{m \times t})}{2m \times t}$$

\odot is Hadamard production and $\sum Q$ stands summing up all elements in Q .

How to learn?

-- Find **W** and **B** that minimize error E

Ideally, find **W** that makes $dE/d\mathbf{W} = 0$ and **B** that makes $dE/d\mathbf{B} = 0$ --- However,

Almost impossible to do --- 1. Multilayer; 2. Not scalar...

Therefore, gradient descending is applied...

$\mathbf{X}_{m \times n}$: Input matrix with m samples and n features

$\mathbf{Y}_{m \times t}$: Target matrix with m samples and t predicted output

\mathbf{Q}^T : Transpose of matrix \mathbf{Q}

L : Output layer

n^a : number of nodes in layer a

$\mathbf{W}_{n^{a-1} \times n^a}^a$: Weight matrix between upper layer $a - 1$ and lower layer a , the dimension of which is $n^{a-1} \times n^a$

$\mathbf{B}_{1 \times n^a}^a$: Bias matrix(vector) between upper layer $a - 1$ and lower layer a , the dimension of which is $1 \times n^a$

$\mathbf{Z}_{m \times n^a}^a$: Weighted input for layer a , the dimension of which is $m \times n^a$

$\mathbf{A}_{m \times n^a}^a$: Output matrix of layer $a - 1$, or, input matrix of layer a

$\sigma^a(\mathbf{Z})$: Activation function of layer a mapping matrix \mathbf{Z} element-wise

E : Error calculated by loss function

$\frac{\partial Y}{\partial X}$: Differential of variable Y on variable X



Gradient descending:

for each iteration:


$$\mathbf{W} = \mathbf{W} - \text{learning_rate} * dE/d\mathbf{W},$$

$$\mathbf{B} = \mathbf{B} - \text{learning_rate} * dE/d\mathbf{B}$$

$$\mathbf{A}_{m \times t}^{L+1} = \sigma^L((\dots \sigma^2(\sigma^1(\mathbf{X}_{m \times n} \cdot \mathbf{W}_{n \times n^1}^1 + \mathbf{B}_{1 \times n^1}^1) \cdot \mathbf{W}_{n^1 \times n^2}^2 + \mathbf{B}_{1 \times n^2}^2) \dots) \dots \mathbf{W}_{n^{L-1} \times t}^L + \mathbf{B}_{1 \times t}^L)$$

$$E = \frac{\sum (\mathbf{A}_{m \times t}^{L+1} - \mathbf{Y}_{m \times t}) \odot (\mathbf{A}_{m \times t}^{L+1} - \mathbf{Y}_{m \times t})}{2m \times t}$$

\odot is Hadamard production and $\sum \mathbf{Q}$ stands summing up all elements in \mathbf{Q} .


$$E = f(\mathbf{A}^{L+1}) = f(g(\mathbf{Z}^L)) = f(g(h(\mathbf{W}^L)))$$

Derivative with Chain rule:

$$d\mathbf{W}^L = \frac{\partial E}{\partial \mathbf{W}} = \frac{\partial E}{\partial \mathbf{A}^{L+1}} \times \frac{\partial \mathbf{A}^{L+1}}{\partial \mathbf{Z}^L} \times \frac{\partial \mathbf{Z}^L}{\partial \mathbf{W}^L} = (\mathbf{A}_{m \times n^L}^{L-1})^T \times \frac{\mathbf{A}_{m \times t}^{L+1} - \mathbf{Y}_{m \times t}}{m \times t} \odot \sigma'^L(\mathbf{Z}^L)$$

$$d\mathbf{B}^L = \frac{\partial E}{\partial \mathbf{B}} = \frac{\partial E}{\partial \mathbf{A}^{L+1}} \times \frac{\partial \mathbf{A}^{L+1}}{\partial \mathbf{Z}^L} \times \frac{\partial \mathbf{Z}^L}{\partial \mathbf{B}^L} = \frac{\mathbf{A}_{m \times t}^{L+1} - \mathbf{Y}_{m \times t}}{m \times t} \odot \sigma'^L(\mathbf{Z}^L) \times 1$$

$$\Delta_{m \times t}^L = d\mathbf{Z}^L = \frac{\partial E}{\partial \mathbf{Z}^L} = \frac{\partial E}{\partial \mathbf{A}^{L+1}} \times \frac{\partial \mathbf{A}^{L+1}}{\partial \mathbf{Z}^L} \times = \frac{\mathbf{A}_{m \times t}^{L+1} - \mathbf{Y}_{m \times t}}{m \times t} \odot \sigma'^L(\mathbf{Z}^L)$$

$$\Delta_{m \times n^l}^l = d\mathbf{Z}^l = \frac{\partial E}{\partial \mathbf{Z}^l} = \Delta^{l+1} \times \frac{\partial \mathbf{Z}^{l+1}}{\partial \mathbf{A}^{L+1}} \times \frac{\partial \mathbf{A}^{L+1}}{\partial \mathbf{Z}^L} = \Delta^{l+1} \cdot (\mathbf{W}_{n^i \times n^{i+1}}^{i+1})^T \odot \sigma'^l(\mathbf{Z}^l)$$

Derivation of E on weighted input

Error on Layer l



4 important equations in back propagation:

$$\Delta_{m \times t}^L = \frac{\mathbf{A}_{m \times t}^{L+1} - \mathbf{Y}_{m \times t}}{m \times t} \odot \sigma'^L(\mathbf{Z}^L)$$

$$\Delta_{m \times n^l}^l = \Delta^{i+1} \cdot (\mathbf{W}_{n^i \times n^{i+1}}^{i+1})^T \odot \sigma'^l(\mathbf{Z}^l)$$

$$d\mathbf{W}^l = (\mathbf{A}^{l-1})^T \times \Delta^l$$

$$d\mathbf{B}^l = \Delta^l$$

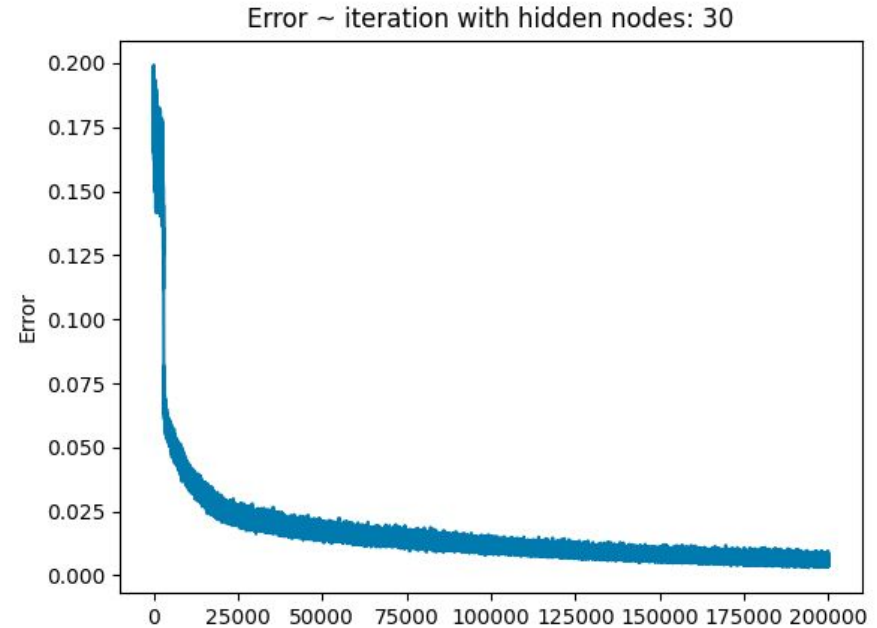
Why back propagation?

Propagate error to each layer and update \mathbf{W} and \mathbf{B} layer by layer backward.



With implemented program,
 $y = (\sin(x_1 + x_2) + 1)/2$
was simulated with 30 hidden nodes
and 10000 samples.

Algorithm can converge with low error
(less than 0.01)!





What to improve?

Activation function, loss function,
structure, training techniques...

Activation functions:

Sigmoid:
range 0 ~ 1;
smooth

tanh:
range -1 ~ 1;
smooth

Relu:
range 0 ~ 1;
Computation
efficient

Softmax:
range 0 ~ 1;
multi-class

Loss functions:

Mean square error:
For real number output

Binary crossentropy:
For binary classification

Categorical Crossentropy:
For multiclass classification

Sparse Categorical Crossentropy:
Same as Category Crossentropy, but
the output formats differ



Chosen datasets:

MNIST -- Image analysis (numbers classification)

Fashion MNIST -- Image analysis (Clothes classification)

CIFAR-10 -- Image analysis (Object detection)

Daily Minimum Temperatures in Melbourne --Time series

**Thank you for your
attention!**