

Neural Network and Backpropagation Questions

Daeyoung Kim, Xintian Han

CDS, NYU

April 21, 2021

[JK] Annotated

Recap: Neural Network and Backpropagation

Feature/representation learning

Initial idea from neuroscience

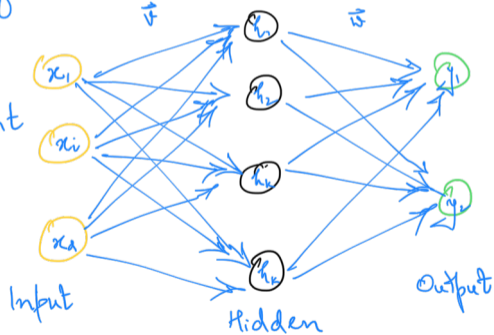
Single hidden layer is sufficient to approximate any fn

Gradient Descent

Special case of automatic differentiation

Algorithm to compute gradient

partial derivatives + chain rule



$$h_k = \sigma(\mathbf{w}_k^T \mathbf{x})$$

$$\begin{aligned} f(\mathbf{x}) &= \sum_{k=1}^K w_k h_k(\mathbf{x}) \\ &= \sum_{k=1}^K w_k \sigma(\mathbf{w}_k^T \mathbf{x}) \end{aligned}$$

Question 1: Step Activation Function ¹

Suppose we have a neural network with one hidden layer.

$$f(x) = w_0 + \sum_i w_i h_i(x); \quad h_i(x) = g(b_i + v_i x),$$

where activation function g is defined as

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

step fn

Which of the following functions can be exactly represented by this neural network?

- polynomials of degree one: $l(x) = ax + b$
- hinge loss: $l(x) = \max(1 - x, 0)$
- polynomials of degree two: $l(x) = ax^2 + bx + c$
- piecewise constant functions

¹From CMU

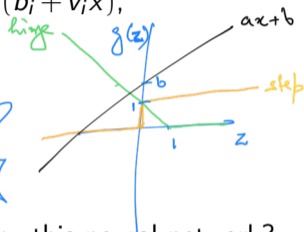
[Solution] Question 1: Step Activation Function

Suppose we have a neural network with one hidden layer.

$$f(x) = w_0 + \sum_i w_i h_i(x); \quad h_i(x) = g(b_i + v_i x),$$

where activation function g is defined as

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$



Which of the following functions can be exactly represented by this neural network?

- polynomials of degree one: $l(x) = ax + b$ **No** — $f(x) = w_0 + \sum_i w_i \mathbb{1}(h_i(x) \geq 0)$
If g can be identity function, then the answer is **Yes** — $f(x) = w_0 + w_1 x + w_2 b$
 $\Rightarrow a = w_1, b = w_0 + w_2 b$
- hinge loss: $l(x) = \max(1 - x, 0)$ **No**
- polynomials of degree two: $l(x) = ax^2 + bx + c$ **No**
- piecewise constant functions **Yes**
 $(-c) \cdot g(x - b) + (c) \cdot g(x - a)$ can represent $l(x) = c, a \leq x < b$.

[Solution] Question 1: Step Activation Function

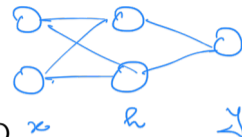
BLANK

Question 2: Power of ReLU ²

Rectified Linear Unit

Consider the following small NN:

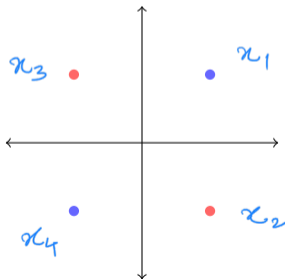
$$w_2^T \text{ReLU}(W_1 x + b_1) + b_2$$



where the data is 2D, W_1 is 2 by 2, b_1 is 2D, w_2 is 2D and b_2 is 1D.

$$x_1 = (1, 1) \quad y_1 = 1; \quad x_2 = (1, -1) \quad y_2 = -1; \quad x_3 = (-1, 1) \quad y_3 = -1; \quad x_4 = (-1, -1) \quad y_4 = 1$$

Find b_1, b_2, W_1, w_2 to solve the problem. (Separate points from class $y = 1$ and $y = -1$.)



²From Harvard

[Solution] Question 2: Power of ReLU

h $f(x)$ $\begin{cases} z & z > 0 \\ 0 & z \leq 0 \end{cases}$

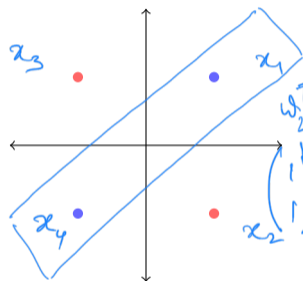
x_1 1 1 $[2 \ 0]$ 1 $w_2^T \text{ReLU}(W_1 x + b_1) + b_2$



x_2 1 -1 $[0 \ 0]$ -1

x_3 -1 1 $[0 \ 0]$ -1

x_4 -1 -1 $[0 \ 2]$ 1



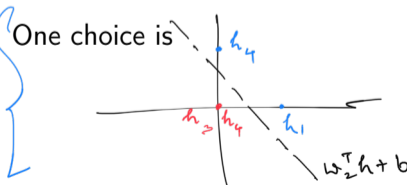
$$\begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix} \begin{bmatrix} x \\ x \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

w_1 x b

$\begin{pmatrix} \pm 1 \\ \pm 1 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \end{pmatrix}$

$+ (-1)$

b_2



$$W_1 = \begin{pmatrix} 1 & 1 \\ -1 & -1 \end{pmatrix}, b_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$w_2 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, b_2 = -1$$

[Solution] Question 2: Power of ReLU

Idea: Map points into different space after performing various operations so as to make them separable / be able to classify

e.g. linear transformations / matrix multiplications

Combination
 $W = U \Sigma V^T$
Singular Value
Decomposition

Rotation: orthonormal matrix
Scaling: diagonal matrix
Reflection: matrix determinant -ve

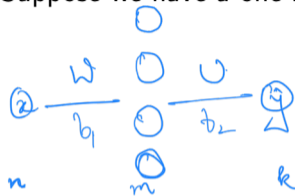
Shearing
Projection
Squeeze

Affine - translation, non-linear transformations (+b)

Reference: <https://atcold.github.io/pytorch-Deep-learning/en/week01/01-3/>

Question 3: Backpropagation ³

Suppose we have a one hidden layer network and computation is:



$$h = \overset{\text{non-linear}}{\text{RELU}}(\overset{\text{affine}}{Wx + b_1})$$

$$\hat{y} = \text{softmax}(Uh + b_2)$$

} 2-layer net

$$J = \text{Cross entropy}(y, \hat{y}) = - \sum_i y_i \log \hat{y}_i$$

n dimension m hidden reps k classes

The dimensions of the matrices are:

$$W \in \mathbb{R}^{m \times n} \quad x \in \mathbb{R}^n \quad b_1 \in \mathbb{R}^m \quad U \in \mathbb{R}^{k \times m} \quad b_2 \in \mathbb{R}^k$$

Use backpropagation to calculate these four gradients

$$\frac{\partial J}{\partial b_2} \quad \frac{\partial J}{\partial U} \quad \frac{\partial J}{\partial b_1} \quad \frac{\partial J}{\partial W}$$

³From Stanford

[Solution] Question 3: Backpropagation

$$\hat{y}_i = \text{softmax}(z_2^{(i)}) = \frac{e^{z_2^{(i)}}}{\sum e^{z_2^{(i)}}}$$

$$\hat{y} = \text{softmax}(z_2)$$

$$\frac{\partial J}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2} \rightarrow \frac{\partial J}{\partial b_2} = \delta_1$$

$$\frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial U} \rightarrow \frac{\partial J}{\partial U} = \delta_1 h^T$$

$$\frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial U} \rightarrow \frac{\partial J}{\partial h} = U^T \delta_1$$

$$z_2 = Uh + b_2 \quad \delta_1 = \frac{\partial J}{\partial z_2} = \hat{y} - y$$

$$\delta_1 = \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2}$$

Cross-entropy = softmax

$$h = \text{ReLU}(z_1)$$

$$\frac{\partial J}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1} \rightarrow \frac{\partial J}{\partial b_1} = \delta_2$$

$$z_1 = Wx + b_1 \quad \delta_2 = \frac{\partial J}{\partial z_1} = U^T \delta_1 \circ 1\{h > 0\}$$

$$\frac{\partial J}{\partial W} = \delta_2 x^T$$

$$\frac{\partial J}{\partial z_1} = \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial h} \cdot \frac{\partial h}{\partial z_1}$$

δ_1 U^T

[Solution] Question 3: Backpropagation

Software

$$\frac{\partial \hat{y}_i}{\partial z_2^{(i)}} = \frac{\partial}{\partial z_2^{(i)}} \left[\frac{e^{z_2^{(i)}}}{\sum_j e^{z_2^{(j)}}} \right]$$

$f_i = j$

$$= \frac{e^{z_2^{(i)}} \sum_j - e^{z_2^{(j)}} e^{z_2^{(i)}}}{\sum_j^2}$$

$$= \frac{e^{z_2^{(i)}}}{\sum_j} - \left(\frac{e^{z_2^{(i)}}}{\sum_j} \right)^2$$

$$\frac{\partial \hat{y}_i}{\partial z_2^{(i)}} = \hat{y}_i - \hat{y}_i^2 = \hat{y}_i (1 - \hat{y}_i)$$

$f_i \neq j$

$$= \frac{0 \sum_j - e^{z_2^{(i)}} e^{z_2^{(j)}}}{\sum_j^2} = - \frac{e^{z_2^{(i)}} e^{z_2^{(j)}}}{\sum_j^2} = - \hat{y}_i \hat{y}_j$$

Quotient rule

for $f(x) = \frac{g(x)}{h(x)}$

$$\rightarrow f'(x) = \frac{g'(x)h(x) - h'(x)g(x)}{[h(x)]^2}$$

Coding Exercise

Cross-Entropy

$$\frac{\partial J}{\partial z_i} = \frac{\partial J}{\partial \hat{y}_i} \left[1 - \sum_j \hat{y}_j \log \hat{y}_j \right]$$

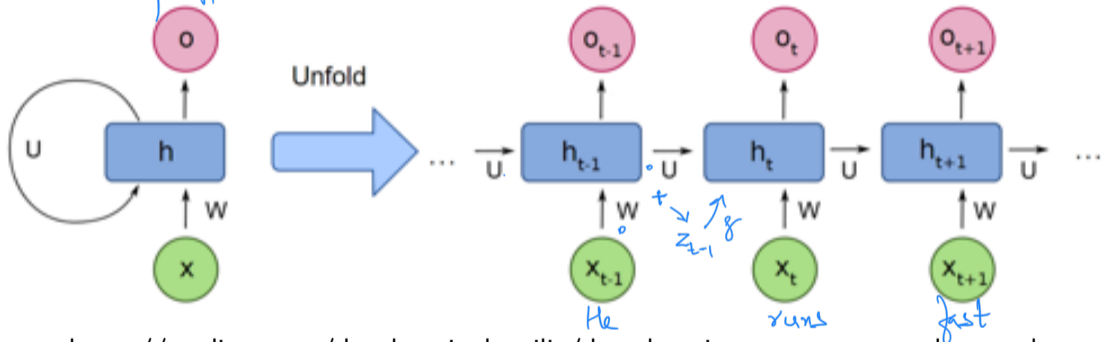
$$= - \frac{\hat{y}_i}{z_i}$$

- Computation graph hands-on

$$\begin{aligned} \frac{\partial J}{\partial z_i^{(i)}} &= \frac{\partial J}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial z_i^{(i)}} = - \frac{\hat{y}_i}{z_i} \frac{\partial \hat{y}_i}{\partial z_i^{(i)}} - \sum_{j \neq i} \frac{\hat{y}_j}{z_i} \frac{\partial (-\hat{y}_j)}{\partial z_i^{(i)}} \\ &= -\hat{y}_i + \hat{y}_i \hat{y}_i - \sum_{j \neq i} \hat{y}_j (-\hat{y}_j) \\ &= \hat{y}_i \left(\hat{y}_i + \sum_{j \neq i} \hat{y}_j \right) = \hat{y}_i \underbrace{\sum_{j=1}^n \hat{y}_j}_{=1} = \hat{y}_i \hat{y}_i = \hat{y}_i^2 \end{aligned}$$

[Optional] Recurrent Neural Networks

For processing variable length i/p, e.g. text-sentences of diff len
Applies recursive fn to each element, e.g. words of sentence



Source: <https://medium.com/deeplearningbrasil/deep-learning-recurrent-neural-networks-f9482a24d010>

Also helps carry over information from prior time steps/position through representations learnt

[Optional]: Backpropagation in RNN

+ BPTT: through time

Suppose we have a recurrent neural network (RNN). The recursive function is:

$$\begin{aligned} \mathbf{z}_{t-1} &= \mathbf{W}\mathbf{x}_{t-1} + \mathbf{U}\mathbf{h}_{t-1}, \\ \mathbf{h}_t &= g(\mathbf{z}_{t-1}), \end{aligned}$$

where \mathbf{h}_t is the hidden state and \mathbf{x}_t is the input at time step t . \mathbf{W} and \mathbf{U} are the weighted matrix. g is an element-wise activation function. And \mathbf{h}_0 is a given fixed initial hidden state.

- Assume loss function \mathcal{L} is a function of \mathbf{h}_T . Given $\partial\mathcal{L}/\partial\mathbf{h}_T$, calculate $\partial\mathcal{L}/\partial\mathbf{U}$ and $\partial\mathcal{L}/\partial\mathbf{W}$.
- Suppose g' is always greater than λ and the smallest singular value of \mathbf{U} is larger than $1/\lambda$. What will happen to the gradient $\partial\mathcal{L}/\partial\mathbf{U}$ and $\partial\mathcal{L}/\partial\mathbf{W}$?
- Suppose g' is always smaller than λ and the largest singular value of \mathbf{U} is smaller than $1/\lambda$. What will happen to the gradient $\partial\mathcal{L}/\partial\mathbf{U}$ and $\partial\mathcal{L}/\partial\mathbf{W}$?

[Solution] [Optional]: Backpropagation in RNN

$$\frac{\partial \mathcal{L}}{\partial U} = \sum_t \frac{\partial \mathcal{L}}{\partial U} ; \quad \frac{\partial \mathcal{L}^{[T]}}{\partial U} = \frac{\partial \mathcal{L}}{\partial h_T} \cdot \frac{\partial h_T}{\partial z_{T-1}} \cdot \frac{\partial (Wx_{T-1} + Uh_{T-1})}{\partial U}$$

$$= \frac{\partial \mathcal{L}}{\partial h_T} \cdot \frac{\partial h_T}{\partial z_{T-1}} \cdot \underbrace{h_{T-1}^T}_{g'(z_{T-1})} + \frac{\partial \mathcal{L}}{\partial h_T} \cdot \frac{\partial h_T}{\partial z_{T-1}} \cdot U^T \cdot \left(\frac{\partial h_{T-1}}{\partial U} \right) + \dots$$

[Solution] [Optional]: Backpropagation in RNN



$$\frac{\partial \mathcal{L}}{\partial U} = \sum_{t=1}^T (\prod_{k=t-1}^{T-1} (\mathbf{U}^T D_k)) \frac{\partial \mathcal{L}}{\partial \mathbf{h}_T} \mathbf{h}_{t-1}^T$$

$$\frac{\partial \mathcal{L}}{\partial W} = \sum_{t=1}^T (\prod_{k=t-1}^{T-1} (\mathbf{U}^T D_k)) \frac{\partial \mathcal{L}}{\partial \mathbf{h}_T} \mathbf{x}_{t-1}^T$$

$D_k = \text{diag}(g'(\mathbf{z}_k))$ is the Jacobian matrix of the element-wise activation function.

- The smallest singular value of the $\mathbf{U}^T D_{k-1}$ will be greater than one. So the smallest singular value of the gradient $\frac{\partial h_s}{\partial h_t}$ will be larger than a^{s-t} for some $a > 1$. So the gradient is going to be exponentially large. This is called **exploding gradient**.
- The largest singular value of the $\mathbf{U}^T D_{k-1}$ will be smaller than one. So the largest singular value of the gradient $\frac{\partial h_s}{\partial h_t}$ will be smaller than a^{s-t} for some $a < 1$. So the gradient is going to be exponentially small. This is called **vanishing gradient**.