

Natural Language Processing with DeepLearning

week 3

NER(Named Entity Recognition)

: 개체명 인식 (자연어 처리 태스크)- 텍스트에서 name을 찾고 분류하는 것

The task: **find** and **classify** names in text, for example:

Last night , Paris Hilton wowed in a sequin gown .

PER PER

Samuel Quinn was arrested in the Hilton Hotel in Paris in April 1989 .

PER PER

LOC LOC LOC DATE DATE

=> 많은 텍스트에서 자동으로 지식 기반 구축을 시작하려는 경우 일반적으로 하고 싶은 것은 명명된 **entity**를 가져와서 이들간의 관계를 파악하는 것

Possible Purpose

- 문서에서 특정 **entity**에 대한 언급 추적
- 질문에 대한 답변의 경우, 답변은 보통 인명
- 요구되는 많은 정보들은 인명간의 관계에 관한 것인 경우가 많음
- 동일한 기술들이 다른 **slot-filling classification**으로 확장될 수 있음

Matrix Calculus

<Jacobian Matrix>

- Given a function with 1 output and n inputs

$$f(\mathbf{x}) = f(x_1, x_2, \dots, x_n)$$

- Given a function with m outputs and n inputs

$$\mathbf{f}(\mathbf{x}) = [f_1(x_1, x_2, \dots, x_n), \dots, f_m(x_1, x_2, \dots, x_n)]$$

$$\frac{\partial f}{\partial \mathbf{x}} = \left[\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]$$

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\left(\frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right)_{ij} = \frac{\partial f_i}{\partial x_j}$$

- 왼쪽(함수 한 개의 미분): n 개의 input을 넣으면 한 개의 output 산출. 이를 미분하면 n 개의 각각의 input으로 미분하여 하나의 벡터가 됨.
- 오른쪽(matrix에서의 미분): m 개의 함수 각각에 n 개의 input이 들어간다고 할 때, 이를 미분하면 input과 output의 모든 조합을 고려하여 $m \times n$ matrix가 되고 이를 Jacobian Matrix라 함

<Chain Rule>

- For composition of one-variable functions: **multiply derivatives**

$$z = 3y$$

$$y = x^2$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx} = (3)(2x) = 6x$$

- For multiple variables at once: **multiply Jacobians**

$$h = f(z)$$

$$z = Wx + b$$

$$\frac{\partial h}{\partial x} = \frac{\partial h}{\partial z} \frac{\partial z}{\partial x} = \dots$$

Stanf

<Example Jacobian: Elementwise activation Function>

Example Jacobian : Elementwise activation Function

$$\begin{aligned} h &= f(z), \text{ what is } \frac{\partial h}{\partial z}? & h, z \in \mathbb{R}^n \\ h_i &= f(z_i) \end{aligned} \quad \longrightarrow \quad \frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Function has n outputs and n inputs $\rightarrow n$ by n Jacobian

H는 z 를 input으로 집어넣어 활성화함수를 적용해 나온 식임. 이를 z 로 미분한다면 우리가 아까 봤던 매트릭스와 비슷한 형태로 $n \times n$ 의 매트릭스 형태가 나옴

$$\begin{aligned} \left(\frac{\partial h}{\partial z} \right)_{ij} &= \frac{\partial h_i}{\partial z_j} = \frac{\partial}{\partial z_j} f(z_i) && \text{definition of Jacobian} \\ &= \begin{cases} f'(z_i) & \text{if } i = j \\ 0 & \text{if otherwise} \end{cases} && \text{regular 1-variable derivative} \end{aligned}$$

$$\frac{\partial h}{\partial z} = \begin{pmatrix} f'(z_1) & & 0 \\ & \ddots & \\ 0 & & f'(z_n) \end{pmatrix} = \text{diag}(f'(z))$$

$z(i)$ 와 $z(j)$ 가 같을 때 미분이 됨
다르면 0으로 없어짐

미분한 식을 알아보면 h_i 는 $f(z_i)$ 로 나타낼 수 있고, z_i 와 z_j 가 $i = j$ 가 같을 때 즉, 변수가 같을 때 미분했을 때 값이 나오게 되고, 다른 경우는 0으로 값이 나오게 됨

이를 행렬로 표현하면 그림과 같이 대각 행렬이 나옴. 이 계산은 우리가 window classification에서 손실 함수를 최소화하기 위한 미분 과정에 부분으로 들어감. 이를 후의 계산에 대입할 것

$$\begin{aligned} \frac{\partial}{\partial x}(Wx + b) &= W \\ \frac{\partial}{\partial b}(Wx + b) &= I \text{ (Identity matrix)} & + & \frac{\partial h}{\partial z} = \begin{pmatrix} f'(z_1) & & 0 \\ & \ddots & \\ 0 & & f'(z_n) \end{pmatrix} = \text{diag}(f'(z)) \\ \frac{\partial}{\partial u}(u^T h) &= h^T \end{aligned}$$

첫 번째는 $wx + b$ 를 x 로 미분한 값 가중치가 나옴, 두 번째는 b 로 미분한 값 항등행렬이 나옴, 세 번째는 u 로 미분한 값 h 를 전치한 값이 나옴.

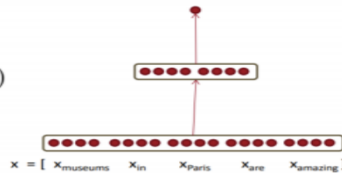
=> 이 세 가지 식과 전에서 봤던 식 모두 우리가 아까 보았던 window classification의 손실 함수를 최소화하기 위해 미분하는 과정에서 사용되는 식들임

Back to our Neural Net!

$$s = \mathbf{u}^T \mathbf{h}$$

$$\mathbf{h} = f(\mathbf{W}\mathbf{x} + \mathbf{b})$$

\mathbf{x} (input)



Let's find $\frac{\partial s}{\partial \mathbf{b}}$

손실함수의 gradient를 계산해야 하지만, 쉽게 score의 gradient를 먼저 계산해보자!

Back to Neural Net! => 손실함수의 gradient를 계산하자 (score함수의 gradient를 계산)
이 score 함수의 gradient를 계산하는 설명은 지금까지의 chain rule과 우리가 계산했던 4가지 식들을 단순히 조합하면 됨.

1. Break up equations into simple pieces

$$s = \mathbf{u}^T \mathbf{h}$$

$$s = \mathbf{u}^T \mathbf{h}$$

$$\mathbf{h} = f(\mathbf{W}\mathbf{x} + \mathbf{b})$$



$$\mathbf{h} = f(\mathbf{z})$$

$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

\mathbf{x} (input)

\mathbf{x} (input)

2. Apply the chain rule

$$\begin{aligned} s &= \mathbf{u}^T \mathbf{h} \\ \mathbf{h} &= f(\mathbf{z}) \\ \mathbf{z} &= \mathbf{W}\mathbf{x} + \mathbf{b} \end{aligned}$$

\mathbf{x} (input)

$$\frac{\partial s}{\partial \mathbf{b}} = \frac{\partial s}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{b}}$$

3. Write out the Jacobians

$$s = \mathbf{u}^T \mathbf{h}$$

$$\mathbf{h} = f(\mathbf{z})$$

$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

\mathbf{x} (input)

$$\begin{aligned} \frac{\partial s}{\partial \mathbf{b}} &= \frac{\partial s}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{b}} \\ &= \mathbf{u}^T \text{diag}(f'(\mathbf{z})) \mathbf{I} \\ &= \mathbf{u}^T \circ f'(\mathbf{z}) \end{aligned}$$

Useful Jacobians from previous slide

$$\begin{aligned} \frac{\partial}{\partial \mathbf{h}} (\mathbf{u}^T \mathbf{h}) &= \mathbf{u}^T \\ \frac{\partial}{\partial \mathbf{z}} (f(\mathbf{z})) &= \text{diag}(f'(\mathbf{z})) \\ \frac{\partial}{\partial \mathbf{b}} (\mathbf{W}\mathbf{x} + \mathbf{b}) &= \mathbf{I} \end{aligned}$$

Re – using Computation

Using the chain rule again:

$$\frac{\partial s}{\partial W} = \frac{\partial s}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial W}$$

$$\frac{\partial s}{\partial b} = \frac{\partial s}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial b}$$

파란색 부분의 계산과정이 같다.
계산을 줄여주는 장점!

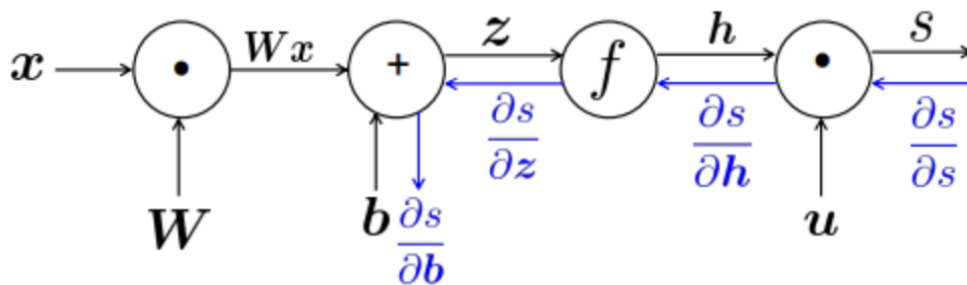
$$\frac{\partial s}{\partial W} = \delta \frac{\partial z}{\partial W}$$

$$\frac{\partial s}{\partial b} = \delta \frac{\partial z}{\partial b} = \delta$$

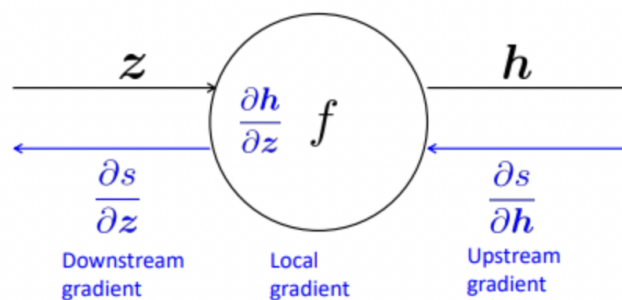
$$\delta = \frac{\partial s}{\partial h} \frac{\partial h}{\partial z} = u^T \circ f'(z)$$

δ is local error signal

Back Propagation



- Backpropagation
Forward Propagation을 통해 얻어진 결과와 실제값을 비교해 계산된 오차를 미분하며 가중치를 업데이트 및 학습.
- Backpropagation 과정에서 중요한 요소
 - (1) Local Gradient
 - (2) Downstream Gradient

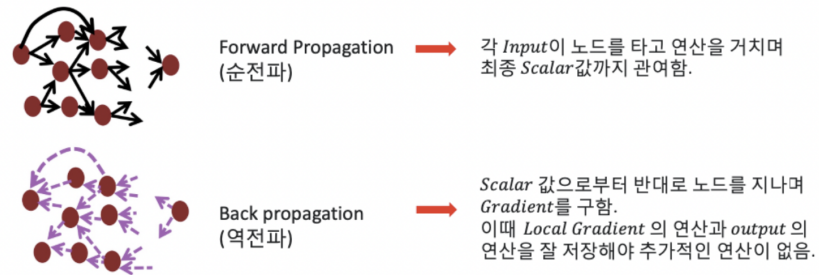


하나의 input과 하나의 output으로 이루어진 single node 그림

Back propagation 과정을 위해 진행 과정에서 Upstream gradient와 Downstream gradient의 순서로 미분값을 계산함

1. Local gradient는 Forward Propagation 수행 시의 output을 input으로 미분
2. 그 다음 Downstream gradient는 Local gradient에 Upstream gradient를 곱하여 도출
결론적으로 Chain Rule을 이용한 것과 같은 결과이다.

<Computation Graph>



검정색 화살표는 순전파가 진행되는 과정

보라색 화살표는 순전파를 통해 나온 최종 값을 바탕으로 오류를 전파하며 거꾸로 값을 업데이트 해가는 역전파과정