



3. Backprop and Neural Networks

1. Named Entity Recognition(NER)

- 텍스트(문장) 내에서 단어를 찾고 분류하는 작업

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

- **context window**를 활용하여 간단히 구현 가능

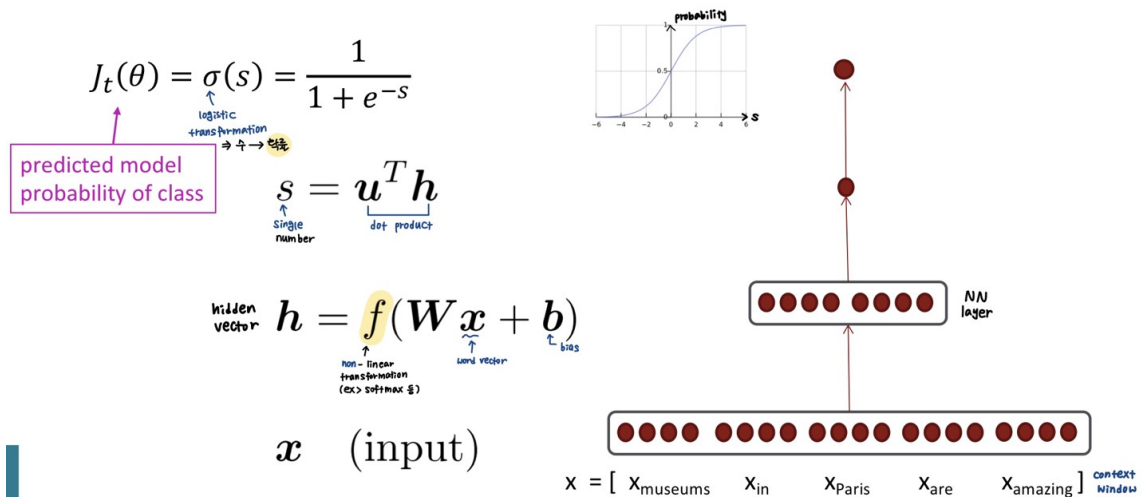
Example: Classify “Paris” as +/- location in context of sentence with window length 2:

the museums in Paris are amazing to see .

$$X_{\text{window}} = [x_{\text{museums}} \quad x_{\text{in}} \quad x_{\text{Paris}} \quad x_{\text{are}} \quad x_{\text{amazing}}]^T$$

NER: Binary classification for center word being location

- We do supervised training and want high score if it's a location



2. Gradients

- input data가 변화 했을 때 output이 얼마만큼 변화하는지를 나타내는 값
 - 각 vector에 대한 편미분 값(partial derivatives)
- multi-input, one-output case
 - 각 input에 대한 partial derivative를 모은 vector가 gradient

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

- multi-input, multi-output case
 - gradient를 matrix 형태로 표현 → Jacobian
 - jacobian matrix는 gradients의 일반화된 형태

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad \boxed{\left(\frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right)_{ij} = \frac{\partial f_i}{\partial x_j}}$$

Chain Rule

- 합성함수의 derivative는 각각의 derivative의 곱으로 얻을 수 있음

- chain rule을 multi-variable 함수에 적용하면 jacobian의 곱으로 확장 가능

$$\begin{aligned} \mathbf{h} &= f(\mathbf{z}) \\ \mathbf{z} &= \mathbf{W}\mathbf{x} + \mathbf{b} \\ \frac{\partial \mathbf{h}}{\partial \mathbf{x}} &= \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \dots \end{aligned}$$

- chain rule을 사용하면 계산을 재사용할 수 있게 됨

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

$$\begin{aligned} \frac{\partial s}{\partial \mathbf{W}} &= \frac{\partial s}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}} \\ \frac{\partial s}{\partial \mathbf{b}} &= \frac{\partial s}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{b}} \end{aligned}$$

The same! Let's avoid duplicated computation ...

Neural Net에 적용하기

1. 각각을 부분으로 나누기

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

2. Chain Rule 적용
3. Jacobian 구하기

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

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

Useful Jacobians from previous slide

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

Derivatives는 어떤 형태여야 할까?

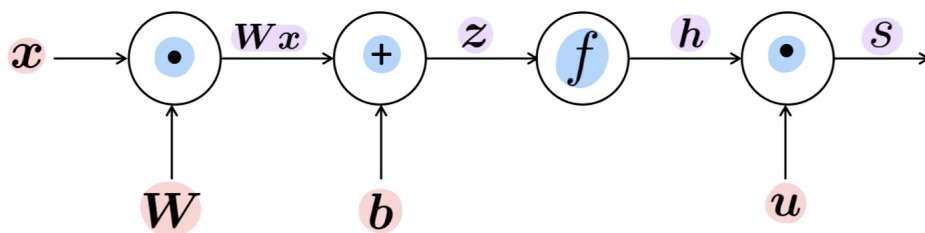
1. Jacobian 형태를 최대한 유지하고, 맨 마지막에 형태 맞추기
2. 형태만을 계속 생각하며 연산하기

3. Backpropagation

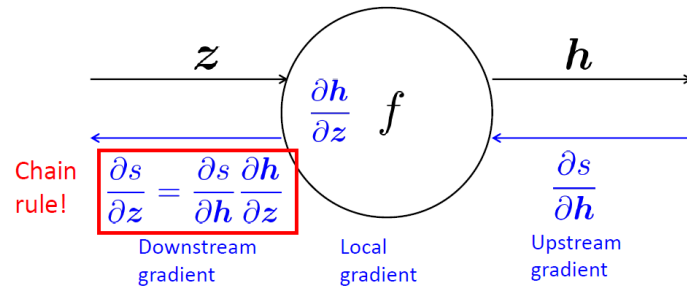
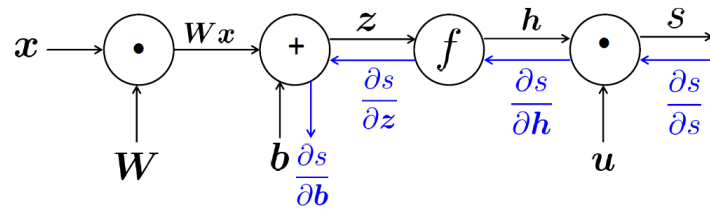
- Neural Network 연산은 아래와 같은 그래프 형태로 표현할 수 있음

- Software represents our neural net equations as a graph
 - Source nodes: inputs
 - Interior nodes: operations
 - Edges pass along result of the operation

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



- 앞에서부터 차례대로 그래프 연산을 진행하는 것을 **forward propagation**이라 함
- output부터 시작하여 역으로 gradients를 계산할 수 있음 ⇒ **back propagation**
 - for parameter update
 - 이때 chain rule을 활용 → 연산된 gradients를 재사용 할 수 있음



- 예시) ppt p.57~