# Named Entity Recognition

- Idea : classify each word in its context window of neighboring words
- dictionary 사용 (×)  NN에 넣어서 특정 label에 대한 확률 계산 (○)
- Ex) "Paris"가 location 인지 window length 2로 살펴보자.

the museums in Paris are amazing to see

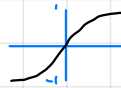$$x_{window} = [\, x_{museums}\ x_{in}\ x_{Paris}\ x_{are}\ x_{amazing}\,]^T$$
$$= x \in \mathbb{R}^{5d}$$

# Non-linearities

- logistic (sigmoid) : $f(z) = \dfrac{1}{1 + \exp(-z)}$

- $\tanh(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}} = 2\,\text{logistic}(2z) - 1$

- $\text{Hard Tanh} = \begin{cases} -1 & \text{if } x < -1 \\ x & \text{if } -1 \leq x \leq 1 \\ 1 & \text{if } x > 1 \end{cases}$

- ReLU (Rectified Linear Unit)
  - Parametric ReLU   GeLU   ReLU

\* Logistic 과 tanh도 쓰이지만 deep networks 에서는 보통 ReLU 사용. GeLU 는 Transformer에 자주 사용

\* non-linearity를 써야 하는 이유 : linear function을 사용하면 layer를 여러개로 하는 의미가 없어짐

# $\nabla_{\theta} J(\theta)$를 계산하는 방법

## 1) By hand (matrix calculus)

### Gradient Computing

n개의 input 이 주어지면 1개의 output이 나오는 함수 $f(x) = f(x_1, \cdots, x_n)$ 에 대하여

$$\frac{df}{dx} = \left[\, \frac{\partial f}{\partial x_1},\ \cdots,\ \frac{\partial f}{\partial x_n}\,\right]$$

m개의 output과 n개의 input 이 주어지는 함수 $f(x) = [\, f_1(x_1, \cdots, x_n),\ \cdots,\ f_m(x_1, \cdots, x_n)\,]$ 에 대하여
(sub function)

$$\frac{df}{dx} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}, \qquad \left(\frac{df}{dx}\right)_{ij} = \frac{\partial f_i}{\partial x_j}$$

$h = f(z)$
$z = W\vec{x} + \vec{b}$ ( $h, z \in \mathbb{R}^n$ ) 처럼 output과 input 의 개수가 동일할 때는 $n \times n$ Jacobian 사용

### Jacobian

$h_i = f(z_i)$

$\left(\dfrac{\partial h}{\partial z}\right)_{ij} = \dfrac{\partial h_i}{\partial z_j} = \dfrac{\partial}{\partial z_j} f(z_i) = \begin{cases} f'(z_i) & \text{if } i = j \\ 0 & \text{if otherwise} \end{cases}$ 정의

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

- other Jacobian

$$\frac{\partial}{\partial x}(W\vec{x} + \vec{b}) = W$$
$$\frac{\partial}{\partial b}(W\vec{x} + \vec{b}) = I$$
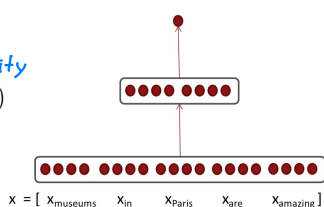$$\frac{\partial}{\partial u}(u^T h) = h^T$$

# Neural Network

$s = u^T h$

non-linearity
$h = f(Wx + b)$

$z = Wx + b$

$x$ (input)

$x = [\, x_{museums}\ x_{in}\ x_{Paris}\ x_{are}\ x_{amazing}\,]$

- loss $J_t$ 의 gradient를 바로 알기 어렵기 때문에 score 의 gradient를 계산

$$\frac{\partial s}{\partial b} = \boxed{\frac{\partial s}{\partial h} \times \frac{\partial h}{\partial z}} \times \frac{\partial z}{\partial b} \quad \Big\}\ \text{by using Jacobians}$$

$$= u^T \text{diag}(f'(z)) \, I$$

$$= u^T \circ f'(z)$$

↳ Hamard Product : element-wise multiplication of 2 vectors to give vector

ex) $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ , $A \odot A = \begin{bmatrix} 1^2 & 2^2 \\ 3^2 & 4^2 \end{bmatrix}$

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

윗 에서와 겹치기 때문에 계산을 줄이기 위해 delta로 바꿈

$$= \delta \frac{\partial z}{\partial b} = \delta \quad (\ \delta = \frac{\partial s}{\partial h} \times \frac{\partial h}{\partial z} = u^T \circ f'(z) \ \text{로 upstream gradient (error signal) 이다.})$$

- 문제점 long row vector로 나타서 계산하기 힘듦

↓ 해결

## Shape Conversion

transpose로 gradient와 parameter들의 shape을 같게 만든다

$$\frac{\partial s}{\partial W} = \delta \frac{\partial z}{\partial W} = \delta \frac{\partial}{\partial W}(Wx + b) = \delta^{\textcircled{T}} x^{\textcircled{T}}$$

↗ upstream gradient
↳ local input signal

# 2) Algorithmically ( Backpropagation )

re-use derivatives computed for higher layers in computing derivatives for lower layers to minimize computation

downstream gradient = upstream gradient × local gradient