# CS-568 Deep Learning

**Nazar Khan**

PUCIT

Automatic Differentiation

## Automatic Differentiation (AD)

▶ Set of techniques to numerically evaluate the derivative of a function *specified by a computer program*.

▶ Analytic or symbolic differentiation evaluates the derivative of a function *specified by a math expression*.

▶ AD Also called *algorithmic differentiation* or *computational differentiation*.

▶ Backpropagation is a special case of AD.

Modern machine learning frameworks (TensorFlow, Theano, PyTorch) employ AD. The programmer only needs to implement the forward pass up to the loss function. Derivatives are handled automatically!

## Automatic Differentiation

*AD exploits the fact that every computer program, no matter how complicated, executes a sequence of elementary arithmetic operations (addition, subtraction, multiplication, division, etc.) and elementary functions (exp, log, sin, cos, etc.). By applying the chain rule repeatedly to these operations, derivatives of arbitrary order can be computed automatically, accurately to working precision, and using at most a small constant factor more arithmetic operations than the original program.*

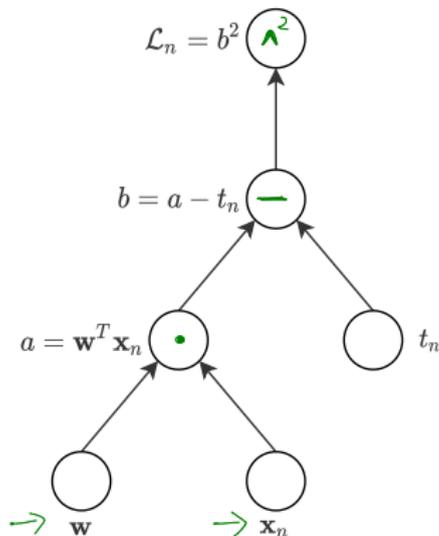*https://en.wikipedia.org/wiki/Automatic_differentiation*

## Linear Regression via Automatic Differentiation

▶ Consider the squared loss function for linear regression.

$$L_{\underline{n}}(\mathbf{w}) = \left(\overbrace{\mathbf{w}^T \mathbf{x}_{\underline{n}}}^{y_n} - t_n\right)^2$$

▶ Can be represented as a computational graph consisting of *elementary operations*.

# Linear Regression via Automatic Differentiation

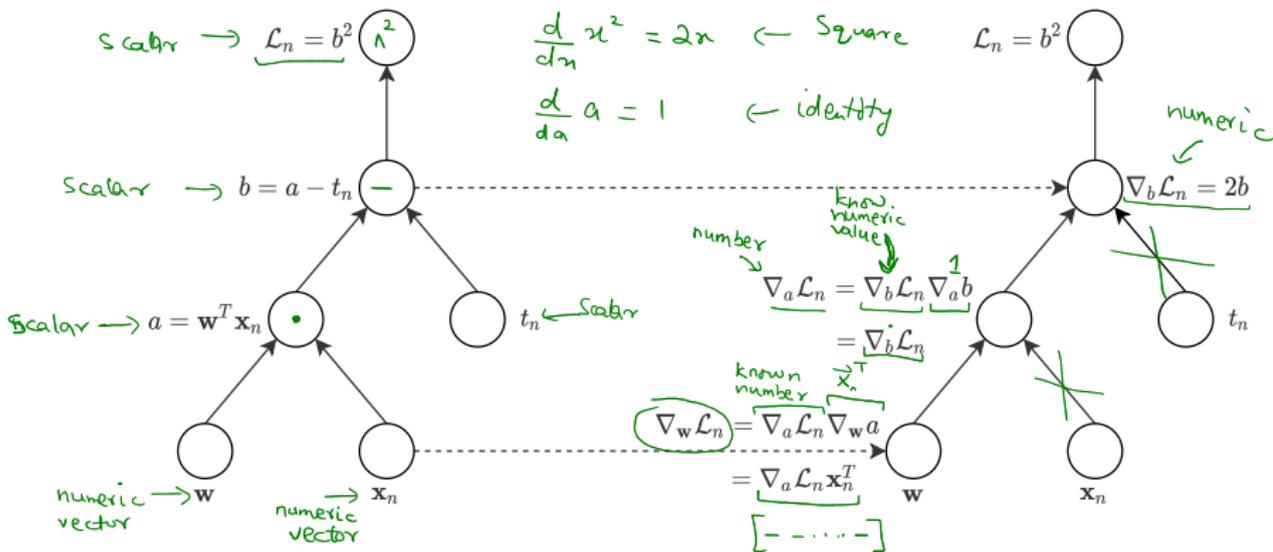$L = (w^T x_n - t_n)^2$    $\boxed{\nabla_w L = 2(w^T x_n - t_n) x_n^T}$    $\left[ \frac{\partial L}{\partial w_1} \; \frac{\partial L}{\partial w_2} \; \cdots \; \frac{\partial L}{\partial w_D} \right]$

<u>No need anymore</u>

▶ For training, we are interested in the gradient $\nabla_{\mathbf{w}} L_n$.

▶ After the forward pass for a particular $\mathbf{w}$ and $\mathbf{x}_n$, gradients can be evaluated numerically.



$\frac{d}{d\vec{w}} \vec{w}^T \vec{x}_n = \vec{x}_n^T$ ← Dot-prod.

$\frac{d}{dn} x^2 = 2x$ ← Square

$\frac{d}{da} a = 1$ ← identity

scalar → $\mathcal{L}_n = b^2$ ⟶ $\boxed{n^2}$       $\mathcal{L}_n = b^2$ ◯

scalar → $b = a - t_n$ ⟶ $\ominus$    known numeric value    ◯ $\nabla_b \mathcal{L}_n = 2b$

numeric

$\nabla_a \mathcal{L}_n = \underline{\nabla_b \mathcal{L}_n} \underline{\nabla_a b}^{\;1}$

scalar → $a = \mathbf{w}^T \mathbf{x}_n$ ⟶ $\odot$    ◯ $t_n$ ← scalar    $= \underline{\nabla_b \mathcal{L}_n}$

◯    ◯ $t_n$

known number $\vec{x}_n^T$

$\boxed{\nabla_{\mathbf{w}} \mathcal{L}_n} = \underline{\nabla_a \mathcal{L}_n} \underline{\nabla_{\mathbf{w}} a}$

numeric → $\mathbf{w}$    $\mathbf{x}_n$    $= \underline{\nabla_a \mathcal{L}_n \mathbf{x}_n^T}$    $\mathbf{w}$    $\mathbf{x}_n$

numeric vector    numeric vector    $[- - \cdots -]$

# AD in Python

▶ A Python package called *Autograd* implements *reverse mode* automatic differentiation.

▶ Elementary operations such as $+, \sin, x^k$ etc. are *overloaded* by also computing their derivates $1, \cos, kx^{k-1}$ etc..

▶ If required, more <u>sophisticated</u> user-defined <u>functions</u> and their <u>derivative</u> <u>implementations</u> can be *registered* with Autograd.

$$\text{elementary functions} \begin{cases} \sigma(a) = \dfrac{1}{1+e^{-a}} \\ \sigma'(a) = \sigma(1-\sigma) \end{cases}$$
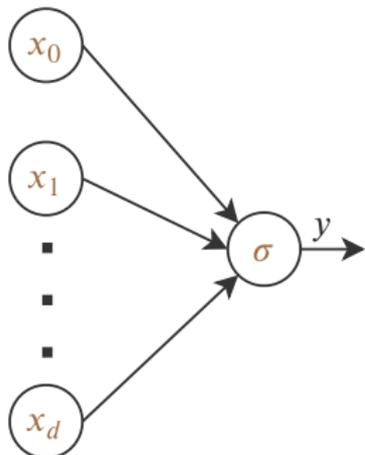
## Logistic Regression via Automatic Differentiation
*Binary classifier with no hidden layer*

Just a perceptron with logistic sigmoid activation function. Models probability of class 1 instead of decision.

$$y = p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T\mathbf{x})$$
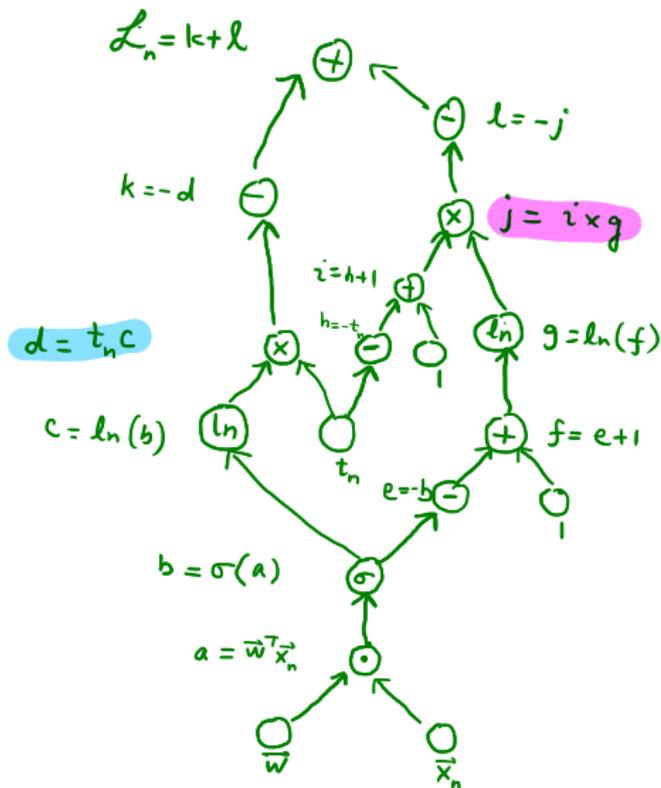$$1 - y = p(\mathcal{C}_2|\mathbf{x}) = 1 - p(\mathcal{C}_1|\mathbf{x})$$



*Binary cross-entropy loss*

$$\mathcal{L}(\mathbf{w}) = -\sum_{n=1}^{N} t_n \ln y_n + (1 - t_n) \ln (1 - y_n)$$

$$\mathcal{L}_n(w) = -t_n \ln y_n - (1-t_n) \ln(1-y_n) \qquad \text{where } y_n = \sigma(\tilde{w}^T \tilde{x}_n)$$

# Logistic Regression via Automatic Differentiation

*Step 1: Computational Graph for $\mathcal{L}_n$*   $\mathcal{L}_n(w) = -t_n \ln y_n - (1-t_n)\ln(1-y_n)$   where $y_n = \sigma(\vec{w}^T \vec{x}_n)$

# Logistic Regression via Automatic Differentiation
*Step 2: AD till $\nabla_{\mathbf{w}} \mathcal{L}_n$*



$$\frac{\nabla}{\mathbf{w}} \mathcal{L}_n \ ?$$