

Vietnam National University, Ho Chi Minh City

University of Science

Faculty of Information Technology

Introduction to Machine Learning

Learning Algorithm

Duc Nguyen

November 2, 2022

Contents

1 Logistic Regression

2 Gradient Descent

3 Automatic Differentiation

Logistic Regression

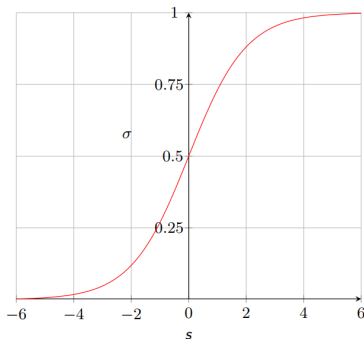
Logistic Function

$$\sigma(s) = \frac{1}{1 + e^{-s}}$$

Properties:

$$\sigma(-s) = 1 - \sigma(s)$$

$$\sigma'(s) = \sigma(s)(1 - \sigma(s))$$



Logistic Regression

Problem statement:

- Objective function f is a probability function

$$f : \mathbb{R}^D \rightarrow [0, 1]$$

- Hypothesis set $h_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$ and probability function

$$P(y|\mathbf{x}, \mathbf{w}) = \begin{cases} h_{\mathbf{w}}(\mathbf{x}) & \text{if } y = 1 \\ 1 - h_{\mathbf{w}}(\mathbf{x}) & \text{if } y = 0 \end{cases}$$

Logistic Regression

Model evaluation

- Likelihood of $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots (\mathbf{x}_n, y_n)\}$:

$$\prod_{n=1}^N P(y_n | \mathbf{x}_n, \mathbf{w})$$

- Maximum a likelihood estimation:

$$\text{Maximize } \prod_{n=1}^N P(y_n | \mathbf{x}_n, \mathbf{w})$$

$$\Leftrightarrow \text{Minimize } -\log \prod_{n=1}^N P(y_n | \mathbf{x}_n, \mathbf{w})$$

Logistic

- Error function:

$$E(h_{\mathbf{w}}) = - \sum_{n=1}^N (y_n \log(h_{\mathbf{w}}(\mathbf{x}_n)) + (1 - y_n) \log(1 - h_{\mathbf{w}}(\mathbf{x}_n)))$$

- **Learning objective:** minimize $E(h_{\mathbf{w}})$
- But how ????

Basic Optimization Problem

$$\begin{aligned} \min_x f(x) \\ \text{subject to } x \in \mathcal{X} \end{aligned}$$

- x is a design point.
- Element in x can be adjusted to minimize the **objective function** f .
- Any value of x from among all points in the **feasible set** \mathcal{X} that minimizes the objective function is called a **solution**.

Local Descent

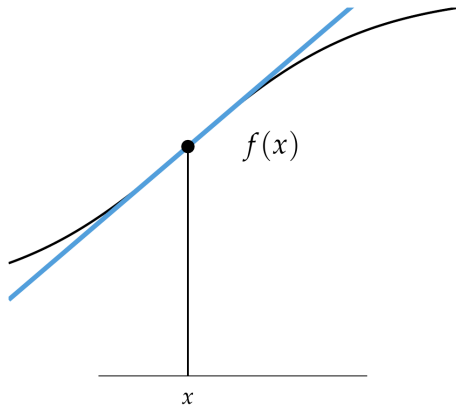
- A common approach for solving an optimization problem is to incrementally improve a design point \mathbf{x} by taking a steps that minimizes the objective value based on a local model
 - 1 Check whether $\mathbf{x}^{(k)}$ satisfies the terminal conditions.
 - 2 Determine the **descent direction** $\mathbf{d}^{(k)}$ using local information.
 - 3 Determine the step size or **learning rate** $\alpha^{(k)}$
 - 4 Compute the next design point:

$$\mathbf{x}^{(k+1)} \leftarrow \mathbf{x}^{(k)} + \alpha^{(k)} \mathbf{d}^{(k)}$$

Gradient Descent

Derivatives I

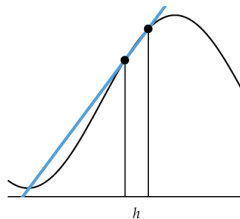
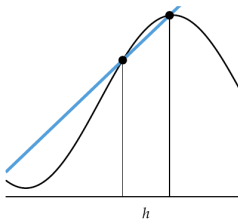
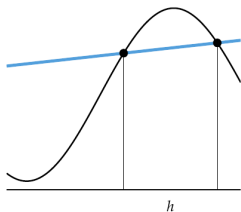
- The derivative $f'(x)$ of a function f of a single variable x is the rate at which the value of f changes at x .
- The value of the derivative equals the slope of the **tangent line**.



Derivatives II

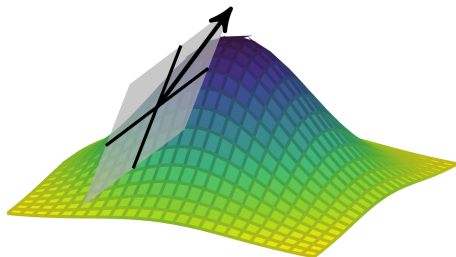
- Derivative can be used provide a linear approximation of the function near x

$$f(x + \Delta x) \approx f(x) + f'(x)\Delta x$$



Derivatives in Multiple Dimensions I

- The **gradient** is the generalization of the derivative to multivariate functions.
- It captures the local slope of the function, allowing us to predict the effect of taking a small step from a point in any direction.



Derivatives in Multiple Dimensions II

- The gradient points in the direction of steepest ascent of the **tangent hyperplane**
- A hyperplane in an n -dimensional space is the set of points that satisfies

$$w_1x_1 + \dots w_nx_n + w_0 = 0$$

$$\mathbf{w}^\top \mathbf{x} = 0$$

- The gradient of f at \mathbf{x} denoted as $\nabla f(\mathbf{x})$ is a vector

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f(\mathbf{x})}{\partial x_0}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right]$$

Gradient Descent

- An intuitive choice for descent direction d is the steepest descent.
- The direction for steepest descent is the direction opposite the gradient ∇f

$$g^{(k)} = \nabla f(x^{(k)})$$

- Typically, we normalize the direction of steepest descent

$$d^{(k)} = -\frac{g^{(k)}}{\|g^{(k)}\|}$$

Logistic Regression (cont.)

Cross-entropy: $J(w) = -(y \log(z) + (1 - y) \log(1 - z))$

Chain rule:
$$\frac{\partial J(w)}{\partial w} = \frac{\partial J(w)}{\partial z} \frac{\partial z}{\partial h} \frac{\partial h}{\partial w}$$

$$\frac{\partial J(w)}{\partial z} = - \left(\frac{y}{z} - \frac{1-y}{1-z} \right) = \frac{z-y}{z(1-z)}$$

$$\frac{\partial z}{\partial h} = z(1-z), \frac{\partial h}{\partial w} = X \rightarrow \frac{\partial J(w)}{\partial w} = X^T(z-y)$$

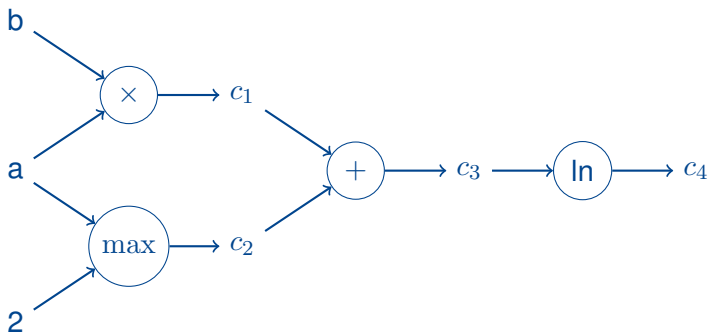
Automatic Differentiation

Automatic Differentiation

- Key to automatic differentiation is the application of Chain rule:

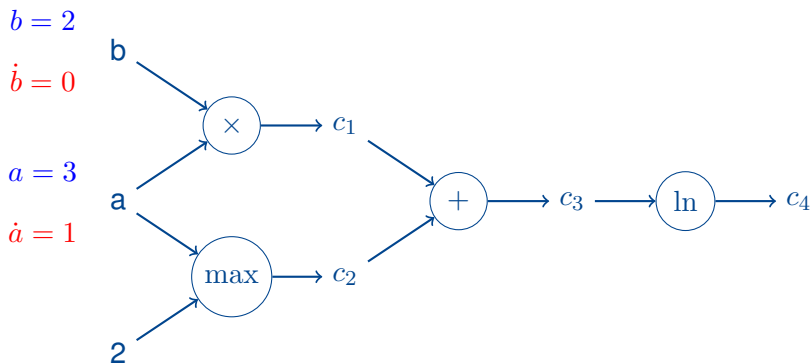
$$\frac{d}{dx} f(g(x)) = \frac{d}{dx} (f \circ g)(x) = \frac{df}{dg} \frac{dg}{dx}$$

Forward Accumulation

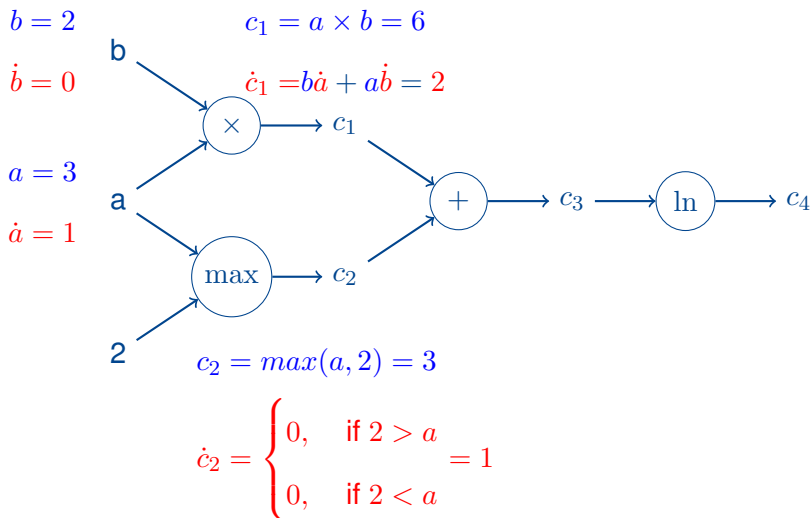


$$\frac{df}{dx} = \frac{df}{dc_4} \frac{dc_4}{dx} = \frac{df}{dc_4} \left(\frac{dc_4}{dc_3} \frac{dc_3}{dx} \right) = \frac{df}{dc_4} \left(\frac{dc_4}{dc_3} \left(\frac{dc_3}{dc_2} \frac{dc_2}{dx} + \frac{dc_3}{dc_1} \frac{dc_1}{dx} \right) \right)$$

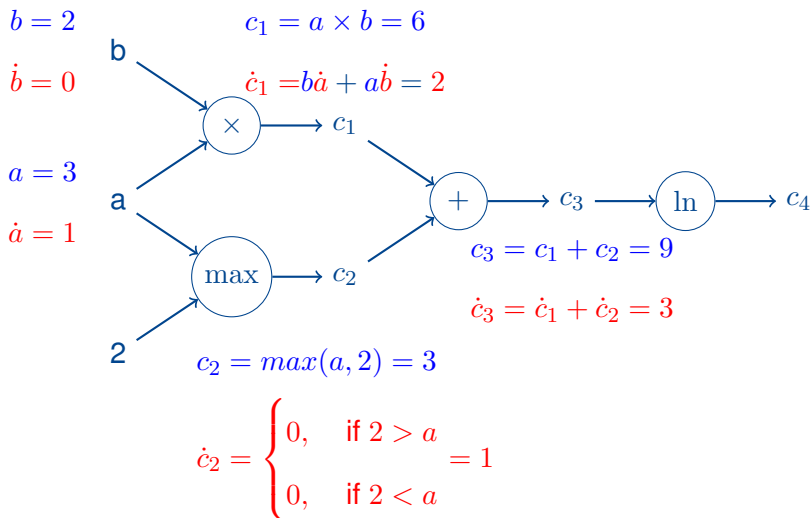
Forward Accumulation



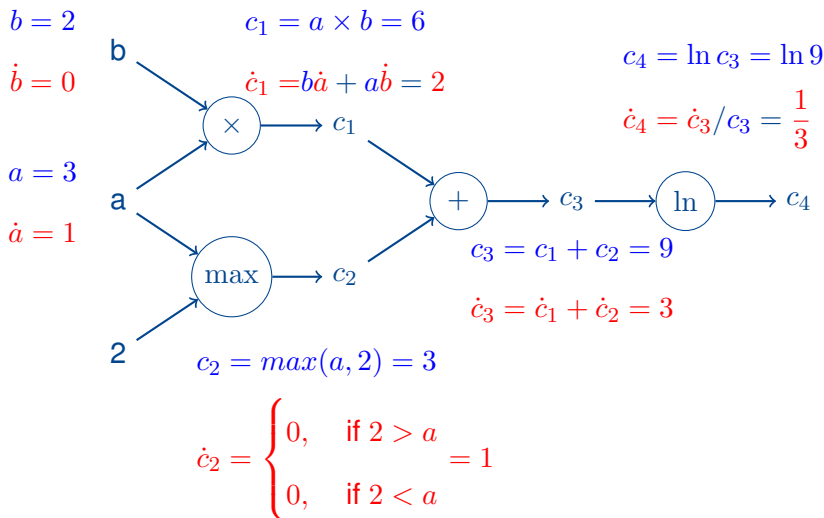
Forward Accumulation



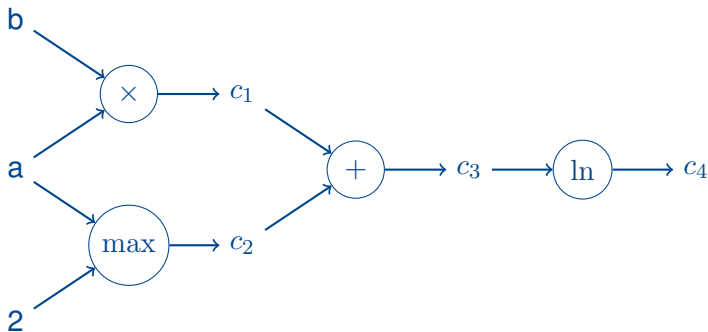
Forward Accumulation



Forward Accumulation

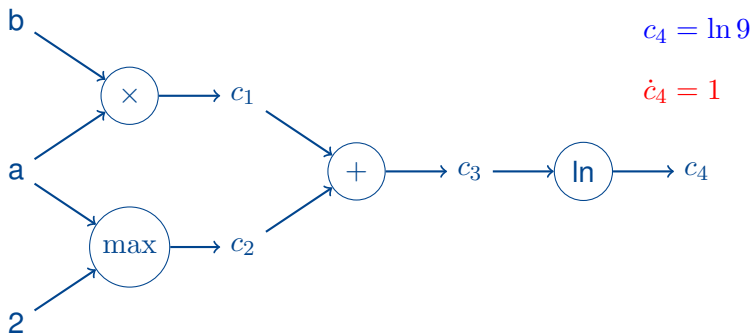


Reverse Accumulation

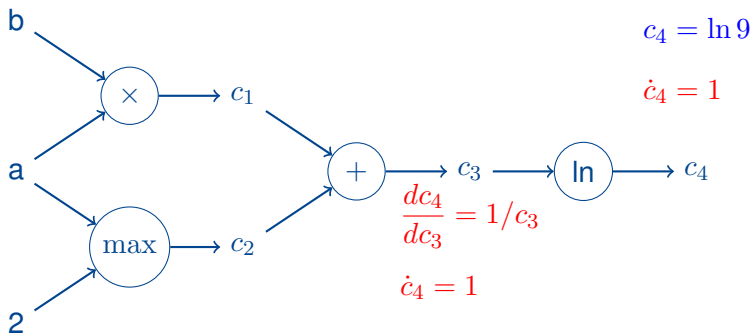


$$\frac{df}{dx} = \frac{df}{dc_4} \frac{dc_4}{dx} = \left(\frac{df}{dc_3} \frac{dc_3}{dc_4} \right) \frac{dc_4}{dx} = \left(\left(\frac{df}{dc_2} \frac{dc_2}{dc_3} + \frac{df}{dc_1} \frac{dc_1}{dc_3} \right) \frac{dc_3}{dc_4} \right) \frac{dc_4}{dx}$$

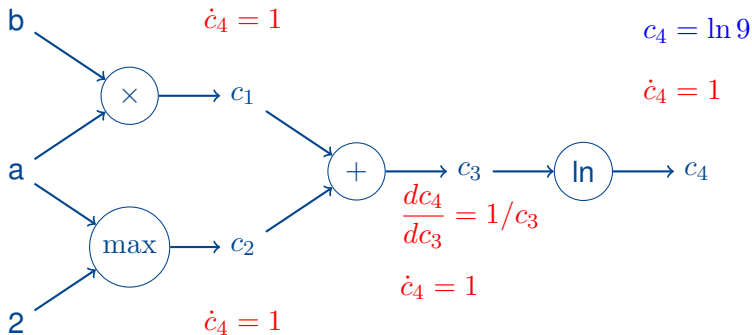
Reverse Accumulation



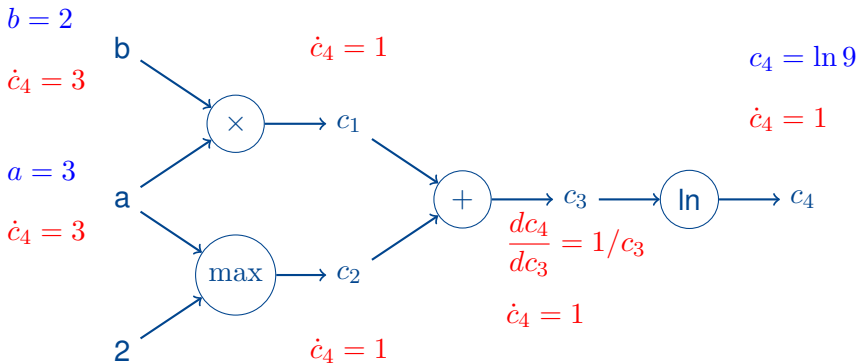
Reverse Accumulation



Reverse Accumulation



Reverse Accumulation



Revision

Revision

- Learning Model: Hypothesis Set, Learning Algorithm.
- Hypothesis Set: **Linear combination** of nonlinear function.
- Learning Algorithm: Gradient Descent.
- Features of training data are the **most important factors** of a learning model.

References



Len Bui.

Course: Introduction to machine learning.
2021.



Gilbert Strang.

Linear algebra and learning from data.
Cambridge Press.