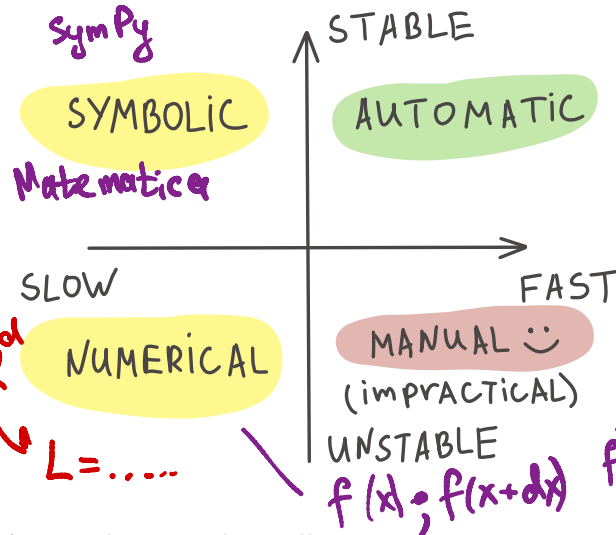
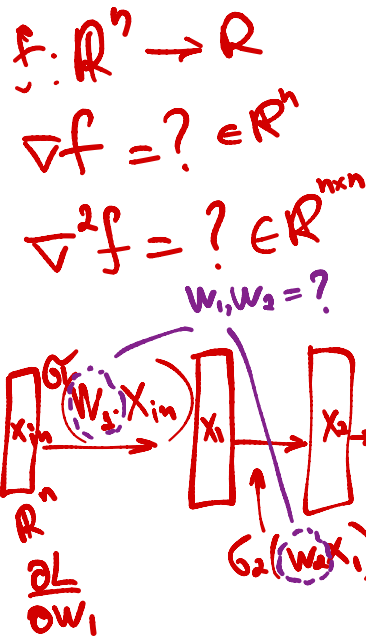


Idea

DIFFERENTIATION



Automatic differentiation is a scheme, that allows you to compute a value of gradient of function with a cost of computing function itself only twice.

Chain rule

We will illustrate some important matrix calculus facts for specific cases

Univariate chain rule

$$W \rightarrow L \rightarrow R$$

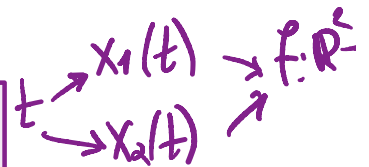
Suppose, we have the following functions $R : \mathbb{R} \rightarrow \mathbb{R}$, $L : \mathbb{R} \rightarrow \mathbb{R}$ and $W \in \mathbb{R}$. Then

$$\frac{\partial R}{\partial W} = \frac{\partial R}{\partial L} \frac{\partial L}{\partial W} \quad R(L(W))$$

Multivariate chain rule

The simplest example:

$$\frac{\partial}{\partial t} f(x_1(t), x_2(t)) = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}$$



Now, we'll consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$:

$$\frac{\partial}{\partial t} f(x_1(t), \dots, x_n(t)) = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \dots + \frac{\partial f}{\partial x_n} \frac{\partial x_n}{\partial t}$$

$$\nabla_x f^T \cdot \frac{dx}{dt}$$

But if we will add another dimension $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, than the j -th output of f will be:

$$\frac{\partial}{\partial t} f_j(x_1(t), \dots, x_n(t)) = \sum_{i=1}^n \frac{\partial f_j}{\partial x_i} \frac{\partial x_i}{\partial t} = \sum_{i=1}^n J_{ji} \frac{\partial x_i}{\partial t},$$

where matrix $J \in \mathbb{R}^{m \times n}$ is the jacobian of the f . Hence, we could write it in a vector way:

$$\frac{\partial f}{\partial t} = J \frac{\partial x}{\partial t} \iff \left(\frac{\partial f}{\partial t}\right)^T = \left(\frac{\partial x}{\partial t}\right)^T J^T$$

$f: \mathbb{R}^n \rightarrow \mathbb{R}^n$
 $f = x + c \quad J = I$

Backpropagation

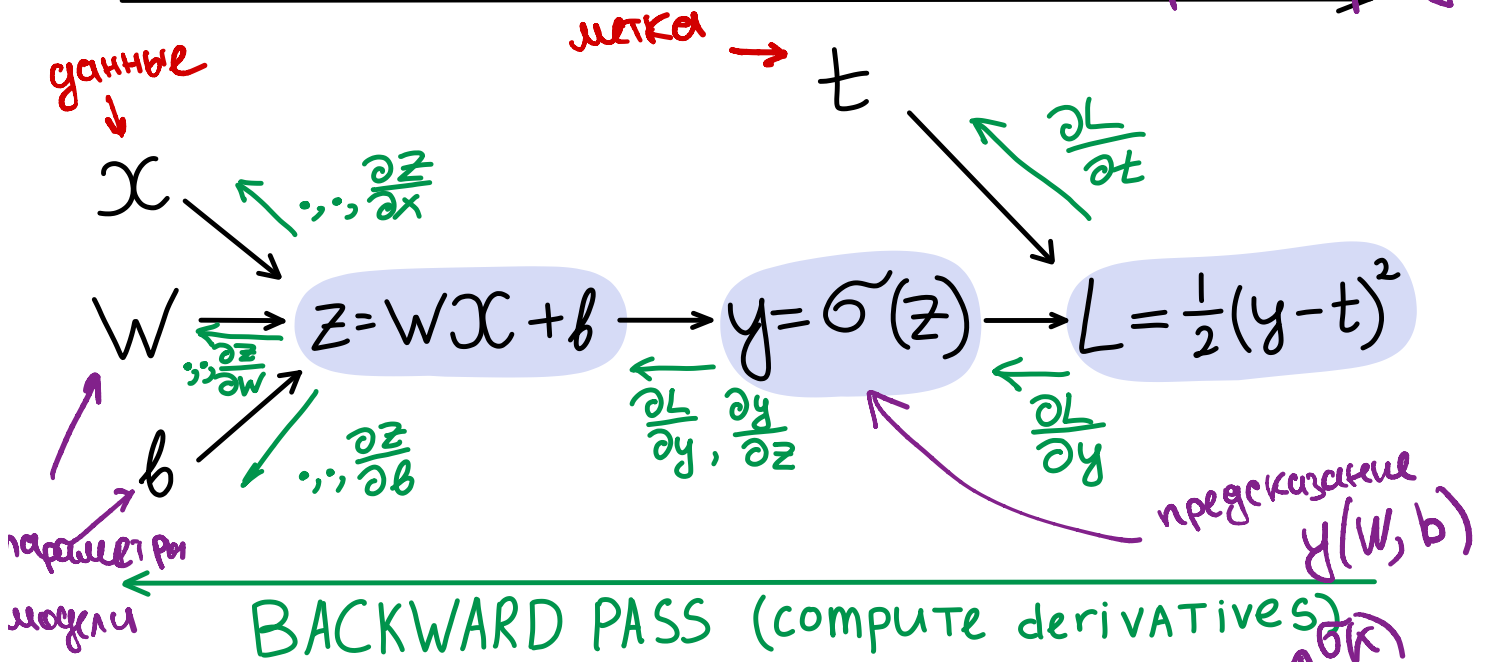
The whole idea came from the applying chain rule to the computation graph of primitive operations

$$L = L(y(z(w, x, b)), t)$$



FORWARD PASS (COMPUTE LOSS)

прямой проход



очу:

$$L = \sum_{k=1}^n W_k - \alpha \frac{\partial L}{\partial W}$$

$$z = wx + b$$

$$\frac{\partial z}{\partial w} = x, \quad \frac{\partial z}{\partial x} = w, \quad \frac{\partial z}{\partial b} = 0$$

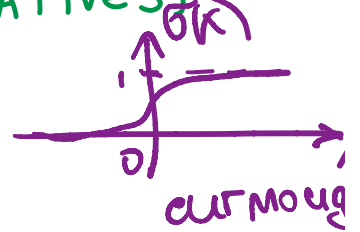
$$y = \sigma(z)$$

$$\frac{\partial y}{\partial z} = \sigma'(z)$$

L

$$L = \frac{1}{2} (y - t)^2$$

$$\frac{\partial L}{\partial y} = y - t, \quad \frac{\partial L}{\partial t} = t - y$$



хотим $y(w, b) \sim f$

All frameworks for automatic differentiation construct (implicitly or explicitly) computation graph. In deep learning we typically want to compute the derivatives of

the loss function L w.r.t. each intermediate parameters in order to tune them via gradient descent. For this purpose it is convenient to use the following notation:

$$\bar{v}_i = \frac{\partial L}{\partial v_i}$$

Let v_1, \dots, v_N be a topological ordering of the computation graph (i.e. parents come before children). v_N denotes the variable we're trying to compute derivatives of (e.g. loss).

Forward pass:

- For $i = 1, \dots, N$:
 - Compute v_i as a function of its parents.

Backward pass:

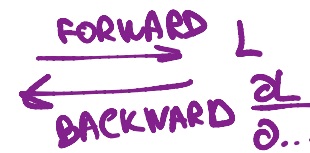
- $\bar{v}_N = 1$
- For $i = N - 1, \dots, 1$:
 - Compute derivatives $\bar{v}_i = \sum_{j \in \text{Children}(v_i)} \bar{v}_j \frac{\partial v_j}{\partial v_i}$

Note, that \bar{v}_j term is coming from the children of \bar{v}_i , while $\frac{\partial v_j}{\partial v_i}$ is already precomputed effectively.

Автомат. уаар.:

1. Строится вычислительный граф примитивных арифм. операц.

2. Для подсчёта производных используются правила прозв. сложной ф(x)

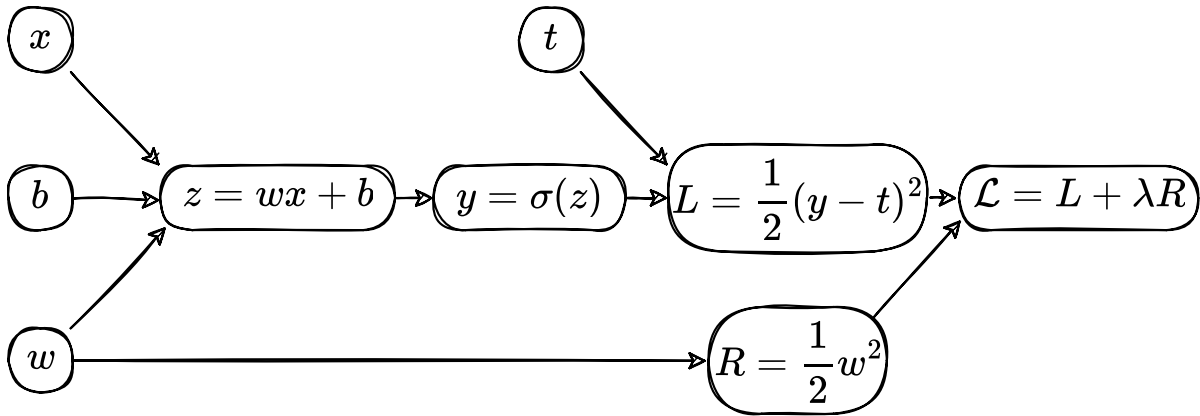


3. Все операции снабжены эффективной реализацией J-вектор JVP (VJP)

FORWARD To

BACKWARD 2.To

Univariate logistic least squares regression



Forward pass

$$\begin{aligned}
 z &= wx + b \\
 y &= \sigma(z) \\
 L &= \frac{1}{2}(y - t)^2 \\
 R &= \frac{1}{2}w^2 \\
 \mathcal{L} &= L + \lambda R
 \end{aligned}$$

Backward pass

$$\begin{aligned}
 \bar{\mathcal{L}} &= 1 \\
 \bar{R} &= \bar{\mathcal{L}} \frac{d\mathcal{L}}{dR} = \bar{\mathcal{L}}\lambda \\
 \bar{L} &= \bar{\mathcal{L}} \frac{d\mathcal{L}}{dL} = \bar{\mathcal{L}} \\
 \bar{y} &= \bar{L} \frac{dL}{dy} = \bar{L}(y - t)
 \end{aligned}$$

$$\begin{aligned}
 \bar{z} &= \bar{y} \frac{dy}{dz} = \bar{y}\sigma'(z) \\
 \bar{w} &= \bar{z} \frac{dz}{dw} + \bar{R} \frac{dR}{dw} = \bar{z}x + \bar{R}w \\
 \bar{b} &= \bar{z} \frac{dz}{db} = \bar{z} \\
 \bar{x} &= \bar{z} \frac{dz}{dx} = \bar{z}w
 \end{aligned}$$

Jacobian vector product

The reason why it works so fast in practice is that the Jacobian of the operations are already developed in effective manner in automatic differentiation frameworks.

Typically, we even do not construct or store the full Jacobian, doing matvec directly instead.

$$\begin{matrix} z_1 \\ \vdots \\ z_n \end{matrix} \Rightarrow \begin{matrix} e^{z_1} \\ \vdots \\ e^{z_n} \end{matrix}$$

Example: element-wise exponent

$$\boxed{y = \exp(z)} \quad \boxed{J = \text{diag}(\exp(z))} \quad \bar{z} = \bar{y}J$$

See the examples of Vector-Jacobian Products from autodidact library:

якобиан умлет
сруктура
=> JVP $\alpha \bar{y}$
VJP

```

defvjp(anp.add,          lambda g, ans, x, y : unbroadcast(x, g),
                                lambda g, ans, x, y : unbroadcast(y, g))
defvjp(anp.multiply,    lambda g, ans, x, y : unbroadcast(x, y * g),
                                lambda g, ans, x, y : unbroadcast(y, x * g))
defvjp(anp.subtract,    lambda g, ans, x, y : unbroadcast(x, g),
                                lambda g, ans, x, y : unbroadcast(y, -g))
defvjp(anp.divide,      lambda g, ans, x, y : unbroadcast(x, g / y),
                                lambda g, ans, x, y : unbroadcast(y, -g * x / y**2))
  
```

```
defvjp(anp.true_divide, lambda g, ans, x, y : unbroadcast(x, g / y),
      lambda g, ans, x, y : unbroadcast(y, - g * x / y**2))
```

Hessian vector product

Interesting, that the similar idea could be used to compute Hessian-vector products, which is essential for second order optimization or conjugate gradient methods. For a scalar-valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with continuous second derivatives (so that the Hessian matrix is symmetric), the Hessian at a point $x \in \mathbb{R}^n$ is written as $\partial^2 f(x)$. A Hessian-vector product function is then able to evaluate

$$v \mapsto \partial^2 f(x) \cdot v$$

for any vector $v \in \mathbb{R}^n$.

The trick is not to instantiate the full Hessian matrix: if n is large, perhaps in the millions or billions in the context of neural networks, then that might be impossible to store. Luckily, `grad` (in the `jax/autograd/pytorch/tensorflow`) already gives us a way to write an efficient Hessian-vector product function. We just have to use the identity

$$\partial^2 f(x)v = \partial[x \mapsto \partial f(x) \cdot v] = \partial g(x),$$

where $g(x) = \partial f(x) \cdot v$ is a new vector-valued function that dots the gradient of f at x with the vector v . Notice that we're only ever differentiating scalar-valued functions of vector-valued arguments, which is exactly where we know `grad` is efficient.

```
import jax.numpy as jnp

def hvp(f, x, v):
    return grad(lambda x: jnp.vdot(grad(f)(x), v))(x)
```

Code

 Open in Colab

Materials

- [Autodidact](#) - a pedagogical implementation of Autograd
- [CSC321 Lecture 6](#)
- [CSC321 Lecture 10](#)
- [Why you should understand backpropagation :\)](#)
- [JAX autodiff cookbook](#)

$$f = -e^{-\langle x, x \rangle} \quad df = \frac{df}{d^2 f} dx$$

$$\text{tr}(AX) = \text{tr}\left[(A^T)^T \cdot X\right] = \langle A^T, X \rangle$$

$$\Rightarrow \nabla f = A^T \quad \cancel{\nabla f = A}$$

$$\text{tr}(I A X) = \langle I, AX \rangle$$

$$df = \langle I, d(AX) \rangle = \langle I, AdX \rangle = \langle A^T, dx \rangle$$