

Matrix calculus

Useful definitions and notations

We will treat all vectors as column vectors by default.

Matrix and vector multiplication

Let A be $m \times n$, and B be $n \times p$, and let the product AB be:

matmul

$$C = AB$$

$$m \times p \quad m \times n \quad n \times p \quad n^2$$

then C is a $m \times p$ matrix, with element (i, j) given by:

$$c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$$



naive

$$O(n^3)$$

Успешно

$$O(n^{\log_2 7})$$

Let A be $m \times n$, and x be $n \times 1$, then the typical element of the product:

matvec

$$z = Ax$$

$$m \times 1 \quad m \times n \quad n \times 1$$

is given by:

$$z_i = \sum_{k=1}^n a_{ik} x_k$$

Finally, just to remind:

$$\bullet C = AB \quad C^T = B^T A^T$$

$$\bullet AB \neq BA$$

$$\bullet e^A = \sum_{k=0}^{\infty} \frac{1}{k!} A^k$$

$$\bullet e^{A+B} \neq e^A e^B$$

$$\bullet \langle x, Ay \rangle = \langle A^T x, y \rangle$$

$$\langle a, b \rangle = \sum_{i=1}^n a_i b_i$$

скал. пр.

Gradient

Gradient Let $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$, then vector, which contains all first order partial derivatives:

$$\nabla f(x) = \frac{df}{dx} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix} \in \mathbb{R}^n$$

Hessian

Let $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$, then matrix, containing all the second order partial derivatives:

Гессуан $\in \mathbb{R}^{n \times n}$

$$f''(x) = \frac{\partial^2 f}{\partial x_i \partial x_j} = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix}$$

But actually, Hessian could be a tensor in such a way: $(f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m)$ is just 3d tensor, every slice is just hessian of corresponding scalar function $(H(f_1(x)), H(f_2(x)), \dots, H(f_m(x)))$

Jacobian

The extension of the gradient of multidimensional $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$:

$$f'(x) = \frac{df}{dx^T} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}$$

Summary

$$f(x) : X \rightarrow Y; \quad \frac{\partial f(x)}{\partial x} \in G$$

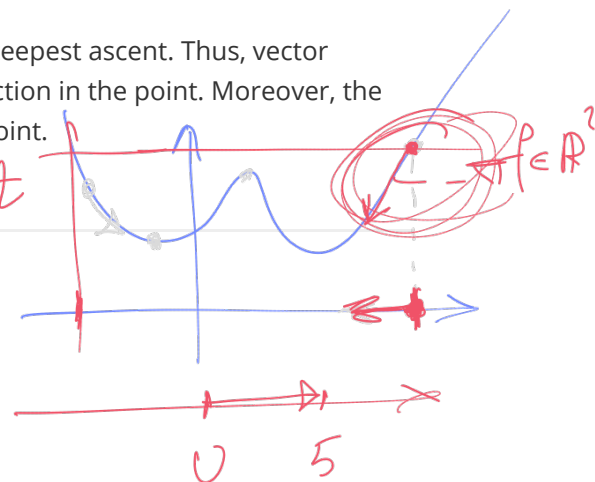
X	Y	G	Name
\mathbb{R}	\mathbb{R}	\mathbb{R}	$f'(x)$ (derivative)
\mathbb{R}^n	\mathbb{R}	\mathbb{R}^n	$\frac{\partial f}{\partial x_i}$ (gradient)
\mathbb{R}^n	\mathbb{R}^m	$\mathbb{R}^{m \times n}$	$\frac{\partial f_i}{\partial x_j}$ (jacobian)
$\mathbb{R}^{m \times n}$	\mathbb{R}	$\mathbb{R}^{m \times n}$	$\frac{\partial f}{\partial x_{ij}}$

named gradient of $f(x)$. This vector indicates the direction of steepest ascent. Thus, vector $-\nabla f(x)$ means the direction of the steepest descent of the function in the point. Moreover, the gradient vector is always orthogonal to the contour line in the point.

General concept

Naive approach

$$f(x) = \text{const}$$



The basic idea of naive approach is to reduce matrix\vector derivatives to the well-known scalar derivatives.

Matrix notation of a function

$$f(x) = c^T x$$

Scalar notation of a function

$$f(x) = \sum_{i=1}^n c_i x_i$$

Matrix notation of a gradient

$$\nabla f(x) = c$$

$$\frac{\partial f(x)}{\partial x_k} = c_k$$

Simple derivative

$$\frac{\partial f(x)}{\partial x_k} = \frac{\partial (\sum_{i=1}^n c_i x_i)}{\partial x_k} = \sum_{i=1}^n \frac{\partial}{\partial x_k} (c_i x_i) =$$

One of the most important practical trick here is to separate indices of sum (i) and partial derivatives (k). Ignoring this simple rule tends to produce mistakes.

$$= \sum_{i=1}^n c_i \cdot \frac{\partial x_i}{\partial x_k} = \sum_{i=1}^n c_i \delta_{ik} = c_k$$

Guru approach

The guru approach implies formulating a set of simple rules, which allows you to calculate derivatives just like in a scalar case. It might be convenient to use the differential notation here.

Differentials

After obtaining the differential notation of df we can retrieve the gradient using following formula:

$$df(x) = \langle \nabla f(x), dx \rangle$$

Then, if we have differential of the above form and we need to calculate the second derivative of the matrix\vector function, we treat "old" dx as the constant dx_1 , then calculate $d(df)$

$$d^2 f(x) = \langle \nabla^2 f(x) dx_1, dx_2 \rangle = \langle H_f(x) dx_1, dx_2 \rangle$$

Properties

Let A and B be the constant matrices, while X and Y are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X+Y) = dX + dY$
- $d(X^T) = (dX)^T$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$
- $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$
- $d(\det X) = \det X \langle X^{-T}, dX \rangle$
- $d\text{tr } X = \langle I, dX \rangle$

$$d(\text{tr } X) = d(\text{tr } (I \cdot X)) = d\langle I, X \rangle$$

$$f(x) = \det X$$

$$f(x+dX) = \det(x+dX)$$

$$df(x) = f(x+dX) - f(x)$$

guru approach

$$\langle A, B \rangle = \text{tr}(A^T B) =$$

$$= \text{tr}(B^T A)$$

- $df(g(x)) = \frac{df}{dg} \cdot dg(x)$
- $H = (J(\nabla f))^T$
- $d(X^{-1}) = -X^{-1}(dX)X^{-1}$

References

- [Good introduction](#)
- [The Matrix Cookbook](#)
- [MSU seminars](#) (Rus.)
- [Online tool](#) for analytic expression of a derivative.
- [Determinant derivative](#)

Example 1

Find $\nabla f(x)$, if $f(x) = \frac{1}{2}x^T A x + b^T x + c$.

Naive approach

$$f(x) = \frac{1}{2} \sum_{i=1}^n x_i (A x)_i + \sum_{i=1}^n b_i x_i + c =$$

$$= \left(\frac{1}{2} \sum_{i=1}^n x_i \left(\sum_{j=1}^n a_{ij} x_j \right) \right) + \sum_{i=1}^n b_i x_i + c$$

$$\frac{\partial f}{\partial x_k} = \frac{1}{2} \frac{\partial}{\partial x_k} \left(\sum_{i,j} x_i a_{ij} x_j \right) + \frac{\partial}{\partial x_k} \left(\sum_{i=1}^n b_i x_i \right) + \frac{\partial}{\partial x_k} c$$

$$\frac{\partial}{\partial x_k} \sum_{i,j} x_i a_{ik} x_k = \sum_i a_{ik} x_i$$

$$\frac{\partial}{\partial x_k} \sum_{i,j} x_k a_{kj} x_j = \sum_j a_{kj} x_j$$

$$2 a_{kk} x_k$$

Example 2

Find $\nabla f(x)$, $f''(x)$, if $f(x) = -e^{-x^T x}$.

$$d(-e^{-\langle x, x \rangle}) = -d(e^{-\langle x, x \rangle}) = -e^{-\langle x, x \rangle} \cdot d(-\langle x, x \rangle) =$$

$$= -e^{-\langle x, x \rangle} \cdot \langle -2x, dx \rangle \Rightarrow \boxed{\nabla f = 2e^{-\langle x, x \rangle} \cdot x}$$

$$\frac{\partial f}{\partial x_k} = \frac{1}{2} \sum_{i=1}^n (a_{ik} x_i + a_{ki} x_i) + b_k =$$

$$= \frac{1}{2} \sum_i (a_{ik} + a_{ki}) x_i + b_k =$$

$$\frac{A+A^T}{2} x + b$$

$$f(x) = \frac{1}{2} \langle x, Ax \rangle + \langle b, x \rangle + c$$

$$df = d\left(\frac{1}{2} \langle x, Ax \rangle + \langle b, x \rangle + c\right) = \frac{1}{2} d(\langle x, Ax \rangle) + d(\langle b, x \rangle) + 0 =$$

$$= \frac{1}{2} \langle dx, Ax \rangle + \frac{1}{2} \langle x, d(Ax) \rangle + \langle b, dx \rangle =$$

$$= \frac{1}{2} \langle Ax, dx \rangle + \frac{1}{2} \langle A^T x, dx \rangle + \langle b, dx \rangle =$$

$$= \langle \frac{1}{2} (A+A^T)x + b, dx \rangle$$

Example 3

Find the gradient $\nabla f(x)$ and hessian $f''(x)$, if $f(x) = \frac{1}{2} \|Ax - b\|_2^2$.

$$\Rightarrow \boxed{\nabla f = 2e^{-\langle x, x \rangle} \cdot x} \Rightarrow df = \langle 2 \cdot e^{-\langle x, x \rangle} \cdot x, dx_1 \rangle = g(x) \quad x_1 = \text{const}$$

$$dg = \langle 2d(e^{-\langle x, x \rangle} \cdot x), dx_1 \rangle = \langle 2(d(e^{-\langle x, x \rangle}) \cdot x + e^{-\langle x, x \rangle} dx), dx_1 \rangle =$$

$$= \langle 2(\langle 2 \cdot e^{-\langle x, x \rangle} \cdot x, dx \rangle \cdot x + e^{-\langle x, x \rangle} dx), dx_1 \rangle =$$

$$= \langle \langle -4tx, dx \rangle x + 2tdx, dx_1 \rangle =$$

$$= \langle x \cdot (-4tx)^T \cdot dx + 2tdx, dx_1 \rangle =$$

$$= \langle -4t \cdot \underbrace{xx^T}_{\mathbb{R}^{n \times n}} dx + 2t \cdot I \cdot dx, dx_1 \rangle \stackrel{!}{=}$$

Example 4

Calculate: $\frac{\partial}{\partial X} \sum \text{eig}(X)$, $\frac{\partial}{\partial X} \prod \text{eig}(X)$, $\frac{\partial}{\partial X} \text{tr}(X)$, $\frac{\partial}{\partial X} \det(X)$

$$\stackrel{!}{=} \langle -2t(2xx^T - I) dx, dx_1 \rangle =$$

$$= \langle (-2t(2xx^T - I))^T dx_1, dx \rangle$$

$$\Rightarrow \boxed{H_f = 2e^{-\langle x, x \rangle} \cdot (I - 2xx^T)}$$

$$X^T X \quad 1 \times n \cdot n \times 1 = 1 \times 1$$

$$X \cdot X^T = n \times n$$

$$XX^T \begin{bmatrix} x_1 & x_2 & \dots & x_n \\ x_1 x_1^T & x_1 x_2^T & \dots & x_1 x_n^T \\ x_2 x_1^T & x_2 x_2^T & \dots & x_2 x_n^T \\ \vdots & \vdots & \ddots & \vdots \\ x_n x_1^T & x_n x_2^T & \dots & x_n x_n^T \end{bmatrix}$$

$$(XX^T)^T = X^T \cdot X = X X^T$$

Example 5

$$X \cdot \rightarrow p \quad X_{\text{new}} = X + \alpha p$$

$$f(X_{\text{new}}) \rightarrow \min_{\alpha=?}$$

$$\nabla_2 f$$

$$\tilde{f} = f(x + \alpha p) \rightarrow \min_{\alpha \in \mathbb{R}} \quad \alpha = ?$$

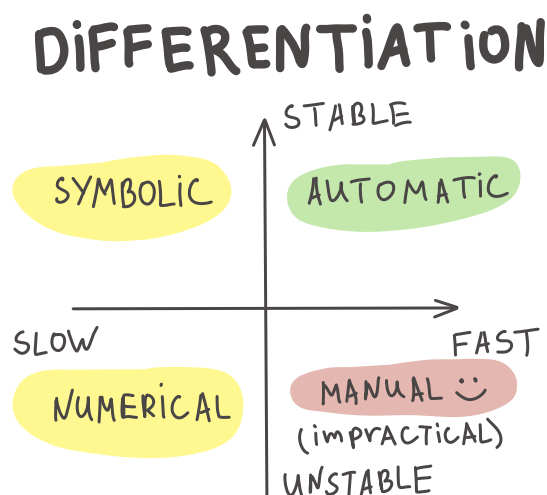
$$df(x + \alpha p) = \langle \nabla_2 f, d\alpha \rangle$$

Find $\nabla f(X)$, if $f(X) = \langle S, X \rangle - \log \det X$

$$\begin{aligned}
 df &= d(\langle S, X \rangle - \ln \det X) = \langle S, dX \rangle - d(\ln \det(X)) = \\
 &= \langle S, dX \rangle - \frac{d(\det X)}{\det X} = \langle S, dX \rangle - \frac{\det X \langle X^{-T}, dX \rangle}{\det X} = \\
 &= \langle S - X^{-T}, dX \rangle \\
 \Rightarrow \boxed{\nabla f = S - X^{-T}} \\
 X^{-T} &= (X^{-1})^T = (X^T)^{-1}
 \end{aligned}$$

Automatic differentiation

Idea



Automatic differentiation is a scheme, that allow 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

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

$$\frac{\partial R}{\partial W} = \frac{\partial R}{\partial L} \frac{\partial L}{\partial W}$$

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 $f : \mathbb{R}^n \rightarrow \mathbb{R}$:

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

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

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

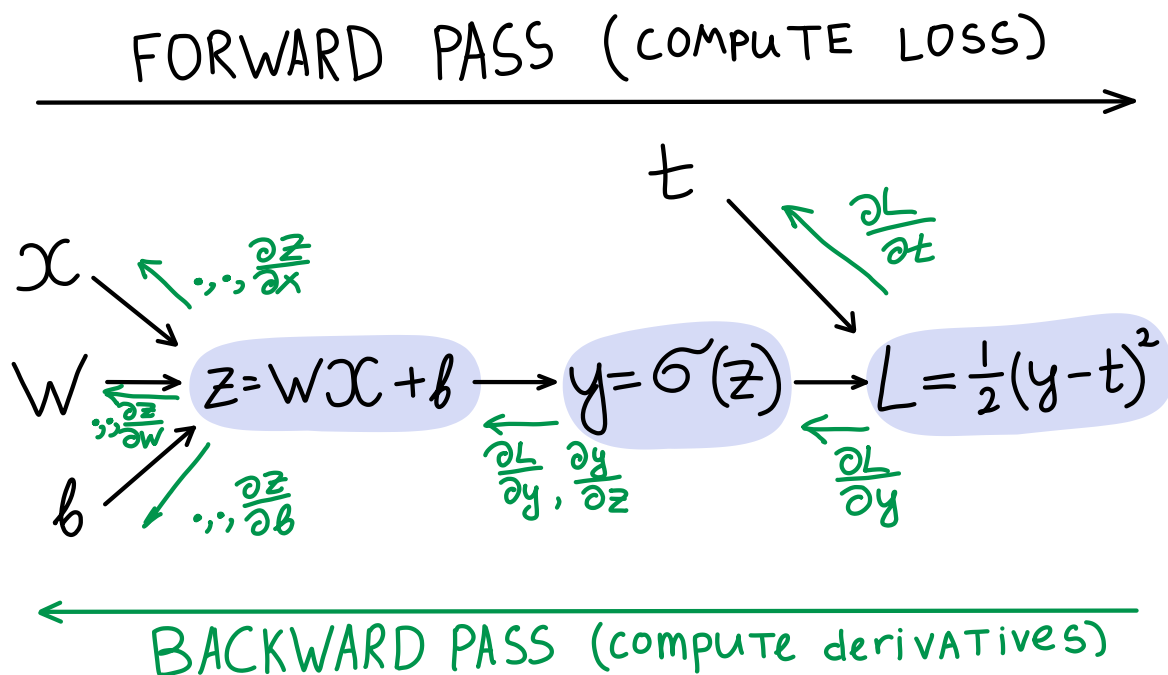
where matrix $J \in \mathbb{R}^{m \times n}$ is the jacobian of the f . Hence, we could write it in a vector way:

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

Backpropagation

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

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



$$\begin{aligned} z &= wx + b & \frac{\partial z}{\partial w} &= x, \frac{\partial z}{\partial x} &= w, \frac{\partial z}{\partial b} &= 0 \\ y &= \sigma(z) & \frac{\partial y}{\partial z} &= \sigma'(z) \\ L &= \frac{1}{2}(y - t)^2 & \frac{\partial L}{\partial y} &= y - t, \frac{\partial L}{\partial t} &= t - y \end{aligned}$$

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:

$$\overline{v_i} = \frac{\partial L}{\partial v_i}$$

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:

- $\overline{v_N} = 1$
- For $i = N - 1, \dots, 1$:
 - Compute derivatives $\overline{v_i} = \sum_{j \in \text{Children}(v_i)} \overline{v_j} \frac{\partial v_j}{\partial v_i}$

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

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.

Example: element-wise exponent

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

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

```
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 scalar-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](#)