

3/20/2023 (1)

Last lecture Chris R showed how to "learn" a differential equation meaning to learn the parameters like learning a neural network.

The derivation was fast & it's worth looking over.

1. Chris's autodiff everything is a vector and a Jacobian

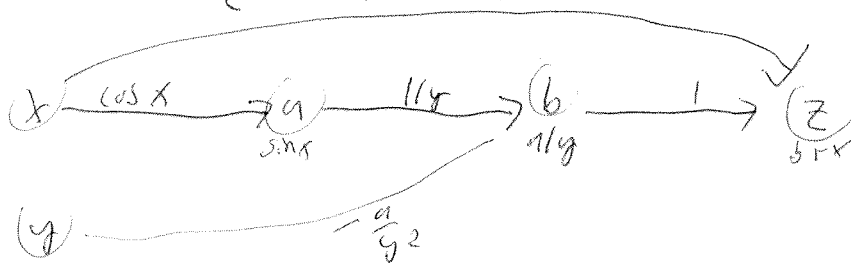
ex1 Graphical as vector-Jacobian

$$\begin{array}{c} (x_1) \\ (x_2) \end{array} \begin{array}{c} \nearrow \frac{\partial y}{\partial x_1} = x_2 \\ \searrow \frac{\partial y}{\partial x_2} = x_1 \end{array} (y) = x_1 x_2$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \xrightarrow{*} y = x_1 x_2$$

$$J = \begin{pmatrix} x_2 & x_1 \end{pmatrix}$$

ex2) x, y $a = \sin x$
 $b = a/y$
 $z = b + x$



Flow written right to left

$$z \leftarrow \begin{pmatrix} b \\ x \end{pmatrix} \leftarrow \begin{pmatrix} a \\ x \\ y \end{pmatrix} \leftarrow \begin{pmatrix} x \\ y \end{pmatrix}$$

$a \neq z$

can be "garbage collected"

3/20/2023 (2)

$$\begin{pmatrix} dx \\ dy \end{pmatrix} = \begin{pmatrix} \cos x & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} dr \\ dg \end{pmatrix}$$

$$\begin{pmatrix} dr \\ dg \end{pmatrix} = \begin{pmatrix} 1/y & 0 & -q/y^2 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} da \\ dx \\ dy \end{pmatrix}$$

$$dz = (1 \quad 1) \begin{pmatrix} dr \\ dg \end{pmatrix}$$

$$dz = (1 \quad 1) \begin{pmatrix} 1/y & 0 & -q/y^2 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} \cos x & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} da \\ dx \\ dy \end{pmatrix}$$

$1 = \nabla_z z$

$$\begin{pmatrix} 1 \\ 1 \end{pmatrix} = (1 \quad 1)^T = \nabla_{\begin{pmatrix} a \\ x \end{pmatrix}} z$$

$$\cancel{(1 \quad 1)^T} \begin{pmatrix} 1/y & 0 \\ 0 & 1 \\ -q/y^2 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \nabla_{\begin{pmatrix} a \\ x \\ y \end{pmatrix}} z$$

$$\begin{pmatrix} \cos x & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1/y & 0 \\ 0 & 1 \\ -q/y^2 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \nabla_{\begin{pmatrix} x \\ y \end{pmatrix}} z$$

(3)

The pullback Function

Let $y = f(x)$ from \mathbb{R}^n to \mathbb{R}^m

If $J \in \mathbb{R}^{m,n}$ is the Jacobian matrix at x

$$J = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{pmatrix}$$

then ~~the pullback~~ $dx \rightarrow J dx$
is known as the push-forward
function. The term comes from
differential geometry and applies to any manifold.

If you have a sequence of functions

$$f_1, f_2, \dots, f_k$$

then

$J_k \cdots J_2 J_1 v$ is the direction
derivative of the composition
in the direction v

and $(J_k \cdots J_1) \begin{pmatrix} y \\ 0 \end{pmatrix}$ is a column of
the total Jacobian

IF $y = f(x)$ (4)

The pullback function

$$\bar{x} = J^T \bar{y} \quad \text{takes } \mathbb{R}^m \text{ to } \mathbb{R}^n$$

It is written $\bar{x} = B_f^x(\bar{y})$

noting that the entries of B_f^x are always functions of x ,

e.g. $f\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} e^{x_1} x_2 \\ \frac{1}{2} \sin(x_1^2 + x_2^2) \\ x_1/x_2 \end{pmatrix}$

$$J = \begin{pmatrix} e^{x_1} x_2 & e^{x_1} \\ \cos(x_1^2 + x_2^2) x_1 & \cos(x_1^2 + x_2^2) x_2 \\ 1/x_2 & -x_1/x_2^2 \end{pmatrix}$$

$$B_f^x(\bar{y}) = \begin{pmatrix} \bar{y}_1 \\ \bar{y}_2 \\ \bar{y}_3 \end{pmatrix} \rightarrow \begin{pmatrix} e^{x_1} x_2 & \cos(x_1^2 + x_2^2) x_1 & 1/x_2 \\ e^{x_1} & \cos(x_1^2 + x_2^2) x_2 & -x_1/x_2^2 \end{pmatrix} \begin{pmatrix} \bar{y}_1 \\ \bar{y}_2 \\ \bar{y}_3 \end{pmatrix}$$

$$\uparrow$$

$$\begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \frac{\partial y_3}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \frac{\partial y_3}{\partial x_2} \end{pmatrix}$$

(5)

As Chris mentioned

$$\begin{pmatrix} J_1^T & \dots & J_{k-1}^T & J_k^T \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \text{ is}$$

a row in the total Jacobian (written as a column)

if f_k is a scalar $J_k^T = \nabla_{x_{k-1}} f_k$

Also if f_k is a scalar f_1, f_2, \dots, f_k is a scalar and

$$J_1^T \dots J_k^T \text{ is } \nabla_{x_1} f_1, \dots, f_k$$

What about matrix to matrix operations like

$Y = X^{-1}$ or matrix to scalar

$$y = \text{tr } M^T X = M \circ X$$

Choice 1:

$$\text{vec}(dY) = - \underbrace{\left(X^{-T} \otimes X^{-1} \right)}_{n^2 \times n^2} \underbrace{\text{vec}(dX)}_{n^2}$$

$J^T = - \left(X^{-1} \otimes X^{-T} \right)$ Kronecker products work that way

$$d\bar{X} = - \left(X^{-1} \otimes X^{-T} \right) \text{vec}(d\bar{Y})$$

(6)

$$X \rightarrow \text{tr } M^T X = M_0 X$$

$$\text{vec}(X) \rightarrow \text{vec}(M) \cdot \text{vec}(X)$$

$$\text{vec}(y) = \text{vec}(M)^T \text{vec}(X)$$

$$\nabla_{\text{vec}(X)} y = \text{vec}(M)$$

Feels Mechanical, which is OK,
I feel adjoints are for grown ups

Recall ~~$\langle u, v \rangle$~~

$\langle u, v \rangle =$ inner product on vector space

$$\langle u, v \rangle = \langle v, u \rangle$$

$$\langle u, u \rangle \geq 0 \quad (= \text{iff } u=0)$$

$$\langle u_1 + u_2, v \rangle = \langle u_1, v \rangle + \langle u_2, v \rangle$$

$$\langle \alpha u, v \rangle = \alpha \langle u, v \rangle \quad \text{et c}$$

Q is a linear operator from V_1 to V_2 with inner products $\langle \cdot, \cdot \rangle_1, \langle \cdot, \cdot \rangle_2$

Q^T is the operator from V_2 to V_1

if ~~$\langle Q^T y_1, y_2 \rangle_1 = \langle y_1, Q y_2 \rangle_2$~~

$$\langle Q^T y, x \rangle_1 = \langle y, Qx \rangle_2 \quad \forall x \in V_1, y \in V_2$$

(7)

e.g. $V_1 = \mathbb{R}^{n,n}$ $V_2 = \mathbb{R}$

$\mathcal{L} = \text{tr } M^T X$ is already linear

$$\langle y_1, y_2 \rangle_{\mathcal{L}} = y_1 y_2$$

$$\langle X_1, X_2 \rangle = X_1^T X_2 = \text{tr } X_1^T X_2$$

$$\begin{aligned} \langle y, \mathcal{L}X \rangle &= y \text{tr } M^T X \\ &= \text{tr } (yM)^T X \\ &= \langle yM, X \rangle \end{aligned}$$

$$\mathcal{L}^T y = My$$

Similarly $\langle X, \mathcal{L}X^T \rangle$

$$\langle Y, -X^{-1} V X \rangle = -\text{tr } Y^T X^{-1} V X$$

$$= -\text{tr } X^T V X^{-T} Y$$

$$= -\text{tr } X^{-1} V X^T Y$$

$$= -\text{tr } X^T V X^{-T} Y$$

$$= -\text{tr } X^{-T} Y X^T V$$

so ~~XAD~~

$$(V \rightarrow -X^{-1} V X)^T$$

=

$$(Y \rightarrow -X^{-T} Y X^T)$$

$$= -\text{tr } ((X^{-T} Y X^T)^T V)$$

(8)

Now $\frac{d}{dx}$ takes functions to functions

$$\langle f, g \rangle = \int_I f(x)g(x) dx \quad \text{is a}$$

common inner product.

Suppose $I=(a,b)$ & $f(a)=f(b)=0$ for simplicity

$$\int_a^b \left[f(x) \frac{d}{dx} g(x) \right] dx = - \int_a^b \left[g(x) \frac{d}{dx} f(x) \right] dx$$

$$\Rightarrow \left(\frac{d}{dx} \right)^T = - \frac{d}{dx}$$

