

We have a sparse matrix 'A' described by a list of 'n' entries 'a' and lists 'ik' and 'jk' of corresponding indices.

To construct A, we define the matrices

$$I = \begin{pmatrix} 1 \\ ik_1, ik_2, \dots \\ 1 \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

... example ...

$$J = \begin{pmatrix} 1 \\ jk_1, jk_2, \dots \\ 1 \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

where the k -th column is a one hot vector into the row/column of A. for entry ' k '.

Then, we may construct A as

$$A_{ij} = I_{ik} J_{jh} a_k$$

or, in matrix notation:

$$A = I \operatorname{diag}(a) J^T$$

where $\operatorname{diag}(a) = \begin{pmatrix} a_1 & & & \\ & a_2 & & \\ & & \ddots & \\ & & & a_n \end{pmatrix}$

for convenience, we define $O = J^T$.

$$A = I \operatorname{diag}(a) O$$

One can confirm that to strip entries 'a_k' from matrix 'A', the process is:

$$a_k = \sum_{ik} \sum_{ij} \sum_{kj} A_{ij} O_{kj}$$

or, in matrix notation:

$$a = \text{diag}(I^T A O^T),$$

where 'diag()' returns the diagonal entries in 'M'.

We have function $b = f(a)$ that takes entries into A and returns entries into $B = AX$.

Our goal is to find the Jacobian

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial a}$$

because $f()$ returns the entries into B , it can be written as

$$f(a) = \text{diag}(I_b^T B O_b^T)$$

$$= \text{diag} \left(I_b^T A X O_b^T \right)$$

we then insert the form of A .

$$= \text{diag} \left(I_b^T I_a \text{diag}(q) O_a X O_b^T \right)$$

The derivative of $\text{diag}(v)$ is just diag of the derivative of v .

$$\text{so } \frac{\partial b}{\partial q} = \text{diag} \left(I_b^T I_a \frac{d \text{diag}(q)}{d q} O_a X O_b^T \right)$$

the term $\frac{d \text{diag}(q)}{d q}$ is a rank 3

tensor with ijk elements $\delta_{ij} \delta_{jk}$.

then, we can write our v_{jps}

as follows.

$$\frac{\partial b}{\partial a}^T v = \text{diag}\left(I_a^T I_b \text{diag}(v) O_b X^T O_a\right)$$

$\text{diag}($

$$\frac{\partial b}{\partial a} g = I_a g D_a X$$

I_b^T

D_b^T

for backprop:

These we need to implement in primitives.

$$\frac{\partial b}{\partial a}^T v = \text{diag}\left(I_a^T \cdot I_b \text{diag}(v) O_b \cdot X^T O_a\right)$$

for -word prop:

$$\frac{\partial b}{\partial a} g = \text{diag}\left(I_b^T \cdot I_a \text{diag}(g) O_a \cdot X \cdot O_b\right).$$

note that $I_c \text{diag}(v) O_c \cdot U$

is equivalent to multiplying
a matrix VU where V uses
entries of V but indices of C .