Matrix Calculus Part II

Linear Transformation View

$$f(x + dx) - f(x) \approx f'(x)[dx]$$

Numerically perhaps take $dx = 1e-7$ or just $.00001$

Abuse of notation:

1. Scalar → Scalar

$$df = \underset{\text{lin transform}}{f'(x)[dx]} = \underset{\text{Scalar}}{f'(x) \cdot dx}$$

e.g. $f(x) = \sin x$    $\{\sin x\}'[dx] = (\cos x) dx$

2. Vector → scalar

$$df = f'(x)[dx] = \underset{\text{row vector}}{f'(x) \cdot dx} \qquad f'(x) = (\nabla_x f)^T$$

e.g. $f(x) = x^T x$

$$\underset{\text{lin transform}}{f'(x)[dx]} \qquad \underset{\text{row vector}}{2x^T dx}$$

3. vector → vector

$$df = f'(x)[dx] = \underset{\text{Jacobian matrix}}{J} dx \qquad (= f'(x) dx) \qquad (J \equiv f_x \equiv \frac{\partial f}{\partial x})$$

4. matrix → Scalar

$$df = f'(x)[dx] = tr\, G^T dx \qquad G = \nabla_x f$$

e.g. $f(x) = tr(x)$    $G = I$

5. matrix → matrix

e.g. $d(x^2)(x^2)'[dx] = x\,dx + dx\,X$

Rules

Plus/Minus

$$F(x) = g(x) \pm h(x)$$
$$f'[dx] = g'[dx] + h'[dx]$$
or $\quad f' = g' + h' \quad$ or $\quad f'(x)[dx] = g'(x)[dx] + h'(x)[dx]$

Product (no matter what kind) $\quad f(x) = g(x) h(x)$
$$f'[dx] = g'[dx] h(x) + g(x) h'[dx]$$

& the most important chain rule
$$F(x) = g(h(x))$$

$$f'(x) = g'(h(x)) \circ h'(x)$$
$$\uparrow$$
composition

or $\quad f'(x)[dx] = g'(h(x)) \left[ h'(x)[dx] \right]$

Try to avoid indices:
the "grown up" approach to matrix calculus

$x \in \mathbb{R}^n:$ $\quad d(x^T x) = x^T dx + dx^T x = 2x^T dx$
dot products of vectors commute

$x \in \mathbb{R}^{n \times n}$ $\quad d(X^2) = X \, dX + dX \, X$
matrices don't commute

Vector to Vector Examples

$x \in \mathbb{R}^2$

1  "rotate by $\theta$"    $f_1(x) = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix} \cdot x$    $\swarrow R(\theta)$

2    warp($\theta$)        $f_2(x) = R(\theta \|x\|) x$
$\uparrow$
the farther from
the origin, the
more you twist

1  $f_1(x) = R(\theta) x$
$d f_1(x) = R(\theta) dx$
$f_1' = R(\theta)$

2    $d\left( R(\theta\|x\|) x \right) = \underline{\underline{\sim}} x + R(\theta\|x\|) dx$

$\sim R'(\theta)$

$d[R(\theta)] = \begin{pmatrix} -\sin\theta & -\cos\theta \\ \cos\theta & -\sin\theta \end{pmatrix} d\theta$

$d\left( (x^T x)^{1/2} \right) = \frac{1}{2} (x^T x)^{-1/2} d(x^T x)$

$= \frac{x^T dx}{\|x\|}$

$d\left( \text{warp}(\theta) \cdot x \right) = \theta R'(\theta\|x\|) x x^T dx + R(\theta) dx$

$= \left[ \theta R'(\theta\|x\|) x x^T + R(\theta) \right] dx$

Chain Rule for vectors

$x \in \mathbb{R}^n$
$\downarrow f$
$y \in \mathbb{R}^K$
$\downarrow g$
$z \in \mathbb{R}^m$
$\downarrow h$
$\ell \in \mathbb{R}$

$$\ell = h(g(f(x)))$$

$$d\ell = J_h(z) \, J_g(y) \, J_f(x) \, dx$$

$$\uparrow \qquad \uparrow \qquad \uparrow$$
$$l \times m \quad m \times K \quad K \times n$$

$$J_h = \nabla_z^T h$$

$$\nabla_x \ell = \nabla_x(h \circ g \circ f) = J_f^T \, J_g^T \, J_h^T$$

$$\uparrow \qquad \uparrow \qquad \uparrow$$
$$n \times K \quad K \times m \quad m \times l$$

Which way is best to form the $m$
numbers of $\nabla_x \ell$ ?

$J_h(J_g J_f)$     forward mode       $2 \left[ (mKn) + (m^2 n) \right]$

$(J_h J_g) J_f$     reverse mode       $2 \left[ mK + mn \right]$

Many inputs one out $\rightarrow$ reverse mode wins
( typical for ML

one input many out $\longrightarrow$ forward mode wins

Notation
   JVP — Jacobian vector product      $f'(x) \cdot v$
   VJP — vector Jacobian product      $v^T f'(x)$

Matrix Linear Operators

$$X \in \mathbb{R}^{n_1, n_2} \qquad Y \in \mathbb{R}^{m_1, m_2}$$

flattening (an always write $\text{vec}(dY) = \boxed{\phantom{xxx}} \text{vec}(dx)$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad m_1 m_2 \times n_1 n_2$

I don't love flattening
but it is kind of a "least common denominator"

Kronecker Product Notation

In Julia: $\quad\quad \text{Kron}(A, B) * \text{vec}(X) \equiv \text{vec}(BXA^T)$

I like $\quad\quad (A \otimes B)[X] = BXA^T$
without writing the vec's

thus $\quad\quad X\,dx = (I \otimes X)[dx]$
$\quad\quad\quad\quad dx\,X = (X^T \otimes I)(dx)$

so $\quad\quad \underset{\text{operator}}{(X^2)'} = \underset{\text{operator}}{\underline{I \otimes X + X^T \otimes I}}$

$$X^{-1} X = I$$
$$d(x^{-1}) X + X^{-1} dx = 0$$
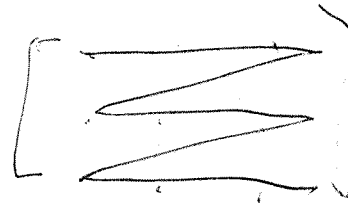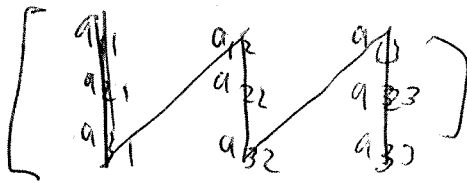$$d(x^{-1}) = -X^{-1} dx\, X^{-1}$$
$$(X^{-1})' = -X^{-T} \otimes X^{-1}$$

④

Part II: Optimizing Serial Code

a)  Column Major    vs   Row Major

Julia                         Python
Matlab                        C
Fortran

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

$$\begin{bmatrix} a_{11} \, a_{21} \, a_{31} & a_{12} \, a_{22} \, a_{32} & a_{13} \, a_{23} \, a_{33} \end{bmatrix} \quad vs \quad \begin{bmatrix} a_{11} \, a_{12} \, a_{13} & a_{21} \, a_{22} \, a_{23} & a_{31} \, a_{32} \, a_{33} \end{bmatrix}$$

✱ Linear Algebra Libraries
  Arrays are linear in Memory

Compare
        for $i=1:n, \quad j=1:n$
  vs  for $j=1:n, \quad i=1:n$
            $c_{ij} = a_{ij} + b_{ij}$
        end

Multicore Memory Architecture
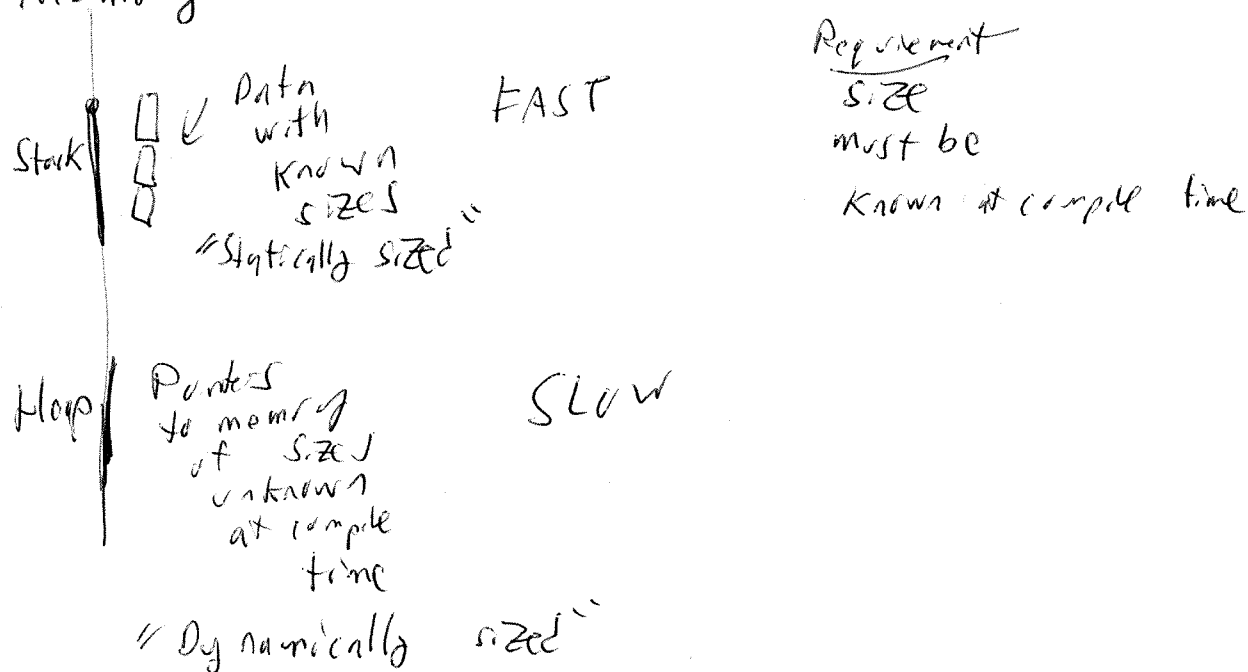    L1, L2 cache

cache miss — when data needs to be
        pulled from main memory

Cache aware algorithm — index/mem structure chosen by
   programmer explicitly to ~~tal~~ avoid cache misses

cache oblivious algor't — indexing structure
   misses cache by design but implicitly


# Main Memory

Stack — Data with Known sizes          FAST

   "Statically sized"

Heap — Pointers to memory of sizes unknown at compile time          SLOW

   "Dynamically sized"

Requirement
Size
must be
Known at compile time

Mutation: use memory already
   preallocated

By convention denoted with a "!" in julia

         Broadcasting
            using  Static Arrays
         @ s   @ SVector

         @ view